

Advanced Statistical Inference

Maurizio Filippone

Lecture 1: Introduction

Recap on linear algebra & probability theory.

Lecture 2: Bayesian Linear Regression

Preliminaries

Probabilities

$$\text{Sum rule: } p(x) = \int p(x, y) dy$$

$$\text{Product rule: } p(x, y) = p(x|y)p(y) = p(y|x)p(x)$$

Expectations

$$\text{Discrete: } \tilde{x} = \mathbb{E}_{p(x)}(x) = \sum xp(x)$$

$$\text{Continuous: } \tilde{x} = \mathbb{E}_{p(x)}(x) = \int xp(x)dx. \text{ In general, } \mathbb{E}_{p(x)}[f(x)] = \int f(x)p(x)dx$$

Mean and Covariance:

$$\mu = \mathbb{E}_{p(x)}[x]$$

$$\sigma^2 = \mathbb{E}_{p(x)}[(x - \mu)^2]$$

$$\text{cov}(x) = \mathbb{E}_{p(x)}[(x - \mu)(x - \mu)^T] = \mathbb{E}_{p(x)}[xx^T] - \mu\mu^T$$

The Gaussian Distribution

$$p(v|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(v - \mu)^2\right\}$$

The Multivariate Gaussian Distribution

$$p(\mathbf{v}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\mathbf{v}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ where } \mathcal{N}(\mathbf{v}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{v} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{v} - \boldsymbol{\mu})\right\}$$

The eigenvalues of the covariance give us the alignment of the distribution.

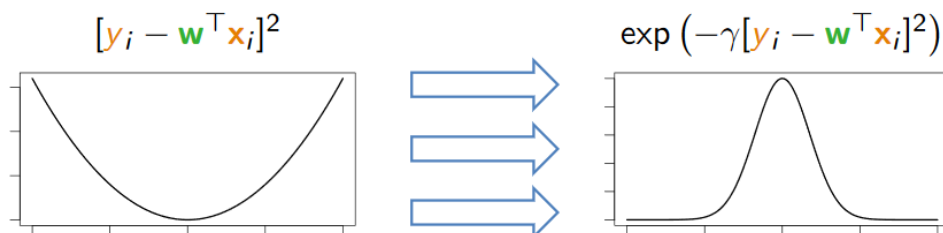
Loss Minimization in Linear Regression

Linear Models for Regression: $f(\mathbf{x}) = \sum w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$

It is a linear model because the parameters appear in a linear way, not because of the linear basis functions.

Quadratic loss function: $\mathcal{L} = \sum (y_i - \mathbf{w}^T \mathbf{x}_i)^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$. Solution: $\nabla_{\mathbf{w}} \mathcal{L} = 0 \Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Probabilistic Interpretation of Loss: minimizing the quadratic loss is equivalent to minimizing the Gaussian likelihood function $\exp(-\gamma \mathcal{L}) = \exp(-\gamma \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2) \propto \mathcal{N}\left(\mathbf{y}|\mathbf{X}\mathbf{w}, \frac{1}{2\gamma}\right)$.



Model Selection

Loss minimization/likelihood maximization is not sufficient to select the best model because we can have generalization problems due to overfitting.

Cross-validation: evaluate models on randomly picked data using validation loss or validation log-likelihood.

Bayesian Inference

Model parameters are considered as probability distributions. Going from $p(\mathbf{w})$ to $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$.

$$\text{Bayes rule: } p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

Posterior density $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$: the distribution over parameters after observing the data.

Likelihood $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$: measure of fitness.

Prior density $p(\mathbf{w})$: anything we know about parameters before we see any data.

Marginal likelihood $p(\mathbf{y}|\mathbf{X})$: normalization constant.

When can we compute the posterior? When the multiplication results is a posterior of same type of density as the prior.

Why it is important? Because we do not have to calculate $p(\mathbf{y}|\mathbf{X})$ because we know the form of $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$.

Finding posterior parameters

The posterior must be Gaussian:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \sigma^2) \propto \exp\left\{-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \boldsymbol{\mu})\right\} \propto \exp\left\{-\frac{1}{2}(\mathbf{w}^T \boldsymbol{\Sigma}^{-1} \mathbf{w} - 2\mathbf{w}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})\right\}$$

On the other hand, we multiply the likelihood and the prior:

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}) &\propto \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})\right\} \exp\left\{-\frac{1}{2}\mathbf{w}^T \mathbf{S}^{-1} \mathbf{w}\right\} \\ &\propto \exp\left\{-\frac{1}{2}\left(\mathbf{w}^T \left[\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \mathbf{S}^{-1}\right] \mathbf{w} - \frac{2}{\sigma^2} \mathbf{w}^T \mathbf{X}^T \mathbf{y}\right)\right\} \end{aligned}$$

We extract parameters:

$$\begin{aligned} \text{Covariance: } \boldsymbol{\Sigma} &= \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \mathbf{S}^{-1}\right)^{-1} \\ \text{Mean: } \boldsymbol{\mu} &= \frac{1}{\sigma^2} \boldsymbol{\Sigma} \mathbf{X}^T \mathbf{y} \end{aligned}$$

Introducing basis functions

We replace \mathbf{X} of size (N, D) by $\boldsymbol{\phi}(\mathbf{X})$ of size (N, D') (expanding or compressing the features).

$$\text{Covariance: } \boldsymbol{\Sigma} = \left(\frac{1}{\sigma^2} \boldsymbol{\phi}^T \boldsymbol{\phi} + \mathbf{S}^{-1}\right)^{-1}$$

$$\text{Mean: } \boldsymbol{\mu} = \frac{1}{\sigma^2} \boldsymbol{\Sigma} \boldsymbol{\phi}^T \mathbf{y}$$

$$\text{Prediction: } p(y_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*, \sigma^2) = \mathcal{N}(y_* | \boldsymbol{\phi}(\mathbf{x}_*)^T \boldsymbol{\mu}, \sigma^2 + \boldsymbol{\phi}(\mathbf{x}_*)^T \boldsymbol{\Sigma} \boldsymbol{\phi}(\mathbf{x}_*))$$

Lecture 3: Gaussian Processes

Introduction

How to choose which basis functions (polynomials, trigonometric, etc.) to use?

Gaussian Processes learn a probabilistic combination of an infinite set of basis functions.

Weight Space View

Recap:

Modeling observations as noisy realizations of a linear combination of features $p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$.

Gaussian prior over model parameters: $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{S})$.

BLR as a Kernel Machine

Working with $\psi(\cdot)$ costs $O(D^2)$ storage and $O(D^3)$ time

Working with $k(\cdot, \cdot)$ costs $O(N^2)$ storage and $O(N^3)$ time

We can pick $k(.,.)$ so that $\psi(.)$ is infinite dimensional.

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x}-\mathbf{x}'\|^2}{2}\right) = \psi(\mathbf{x})^T \psi(\mathbf{x}') \text{ where } \psi(\mathbf{x}) = \exp\left(-\frac{\mathbf{x}^2}{2}\right) \left(1, \mathbf{x}, \frac{\mathbf{x}^2}{\sqrt{2!}}, \dots\right)^T$$

When working with a kernel, we are implicitly working with a finite number of basis functions.

Using Woodbury identity, we can write:

$$\Sigma = \left(\frac{1}{\sigma^2} \Phi^T \Phi + \mathbf{S}^{-1}\right)^{-1} = \mathbf{S} - \mathbf{S} \Phi^T (\sigma^2 \mathbf{I} + \Phi \mathbf{S} \Phi^T)^{-1} \Phi \mathbf{S}$$

We can rewrite the variance:

$$\sigma^2 + \phi_*^T \Sigma \phi_* = \sigma^2 + \phi_*^T \mathbf{S} \phi_* - \phi_*^T \mathbf{S} \Phi^T (\sigma^2 \mathbf{I} + \Phi \mathbf{S} \Phi^T)^{-1} \Phi \mathbf{S} \phi_* = \sigma^2 + k_{**} - \mathbf{k}_*^T (\sigma \mathbf{I} + \mathbf{K})^{-1} \mathbf{k}_*$$

We can also rewrite the mean by applying Woodbury identity twice (no thanks!).

$$\phi_*^T \mu = \dots = \phi_*^T \mathbf{S} \Phi^T (\sigma^2 \mathbf{I} + \Phi \mathbf{S} \Phi^T)^{-1} \mathbf{t} = \mathbf{k}_*^T (\sigma \mathbf{I} + \mathbf{K})^{-1} \mathbf{t}$$

Where

$$\begin{aligned} \psi(\mathbf{x}) &= \mathbf{S}^{1/2} \phi(\mathbf{x}) \\ k_{**} &= k(\mathbf{x}_*, \mathbf{x}_*) = \psi(\mathbf{x}_*)^T \psi(\mathbf{x}_*) \\ (\mathbf{k}_*)_i &= k(\mathbf{x}_*, \mathbf{x}_i) = \psi(\mathbf{x}_*)^T \psi(\mathbf{x}_i) \\ (\mathbf{K})_{ij} &= k(\mathbf{x}_i, \mathbf{x}_j) = \psi(\mathbf{x}_i)^T \psi(\mathbf{x}_j) \end{aligned}$$

ϕ is a D by infinite matrix.

Function Space View

We will consider an infinite number of random variables $f(\mathbf{x})$ indexed by \mathbf{x} . The covariance of those variables is an infinite by infinite zero matrix except on the diagonal where we have the variance of each variable.

$$k(\mathbf{x}, \mathbf{x}') = \alpha \exp(-\beta \|\mathbf{x} - \mathbf{x}'\|^2)$$

Using the distance kernel we can impose nearby variables to have high covariance.

And then we select N variables by selecting the corresponding rows and columns from the covariance matrix.

Marginal distribution $\mathbf{f} = (f(x_1), \dots, f(x_N))^T$ is $p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K})$

Conditional distribution of f_* given \mathbf{f} is $p(f_*|\mathbf{f}, \mathbf{x}_*, \mathbf{X}) = \mathcal{N}(\bar{m}, \bar{s}^2)$

$$\begin{aligned} \text{Where} \quad \bar{m} &= \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{f} \\ \bar{s}^2 &= k_{**} - \mathbf{k}_*^T \mathbf{K}^{-1} \mathbf{k}_* \end{aligned}$$

(posterior and prediction part is missing)

Lecture 4: Bayesian Classification

Classification algorithms:

- Probabilistic: Bayes classifier, Logistic regression.
- Non-probabilistic: K-nearest neighbors, Support Vector Machines.
- Many others ...

Probabilistic classifiers produce a probability of class membership: $P(t_{new} = k | x_{new}, X, t)$

Non-probabilistic classifiers produce a hard assignment, e.g. $t_{new} = 1$ or $t_{new} = 0$.

Probabilities provide us with more information: $P(t_{new} = 1) = 0.6$ is more useful than $t_{new} = 1$.

This is very important when misclassification cost is high and imbalanced (diseased/healthy person).

Logistic Regression

We can apply a logistic function $h(.)$ To $f(x_{new}; \mathbf{w}) = \mathbf{w}^T x_{new}$ to squash it between 0 and 1.

$$\text{We get } P(t_{new} = k | x_{new}, X, t) = h(f(x_{new}; \mathbf{w})) = h(\mathbf{w}^T x_{new}) = \frac{1}{1 + \exp(-\mathbf{w}^T x_{new})}$$

Defining a prior: $p(\mathbf{w}) = \prod_{d=1}^D \mathcal{N}(0, \sigma^2)$

Defining a likelihood:

Assume independence: the noises of observations are independent: $p(t, X, w) = \prod p(t_n | x_n, w)$

For a binary classification: $P(t_n = 0 | x_n, w) = 1 - P(t_n = 1 | x_n, w) = 1 - h(w^T x_n)$

Posterior: $p(w | X, t, \sigma^2) = \frac{p(t | X, w)p(w)}{p(t | X)} = \frac{p(t | X, w)p(w)}{\int p(t | X, w)p(w)dw}$ (can't calculate the integral $\int \left(\frac{1}{1+\exp}\right)^n \dots$)

We can compute $p(t | X, w)p(w) = g(w, X, t)$

We have then three options:

1. Find the most likely value of w
2. Approximate $p(w | X, t)$
3. Sample from $p(w | X, t)$

MAP estimate

$g(w, X, t) \propto p(w | X, t)$. So, if \hat{w} maximizes g it also maximizes the prior (gradient descent).

Once we have \hat{w} , we make prediction with $h(\hat{w}^T x_{new})$

We get a linear boundary (linear model): $h(w^T x_{new}) = 0.5 \Leftrightarrow w^T x_{new} = 0$.

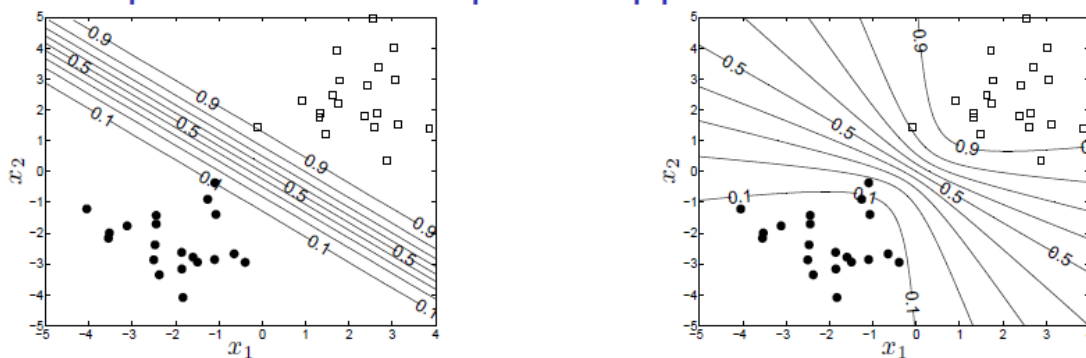
We can approximate $p(w | X, t)$ with a Gaussian $q(w | X, t) = \mathcal{N}(\mu, \Sigma)$, where:

- $\mu = \hat{w} = \arg \max \log g$ ($\arg \max p$ should be a good approximation of the Gaussian mean)
- $\Sigma^{-1} = -\frac{\partial^2 \log g}{\partial w \partial w^T}$ (if g is Gaussian, we get directly the covariance matrix)

Prediction: $P(t_{new} = 1 | x_{new}, X, t) = \mathbb{E}_{\mathcal{N}(\mu, \Sigma)}(P(t_{new} = 1 | x_{new}, w)) = \int \mathcal{N}(\mu, \Sigma) \frac{1}{1+\exp(-w^T x_{new})} dw$ (still can't calculate this integral).

The solution is to sample (we know how to sample from Gaussian) and then average.

Draw S samples $\{w_1, \dots, w_S\}$ from the distribution: $\mathbb{E}_{\mathcal{N}(\mu, \Sigma)}(P(t_{new} = 1 | x_{new}, w)) \approx \frac{1}{S} \sum \frac{1}{1+\exp(-w_s^T x_{new})}$



Bayesian Classifier

Based on Bayes rule: $P(t_{new} = k | X, t, x_{new}) = \frac{P(x_{new} | t_{new} = k, X, t)P(t_{new} = k)}{\sum P(x_{new} | t_{new} = j, X, t)P(t_{new} = j)}$

- Likelihood $P(x_{new} | t_{new} = k, X, t)$: how likely to observe x_{new} if it's in class k ? We can choose this distribution as we like depending on our data (Gaussian, binomial likelihood) \rightarrow Training data with $t = k$ is used to determine params of likelihood for class k (e.g. mean and covariance).
- Prior $P(t_{new} = k)$: this is not related to data (there is no x_{new}).
 - There are fewer instances of class $k = 0$ means $P(t_{new} = 0) < P(t_{new} = 1)$.
 - No prior preference means $P(t_{new} = 0) = P(t_{new} = 1)$.
 - Class 0 is very rare means $P(t_{new} = 0) \ll P(t_{new} = 1)$.

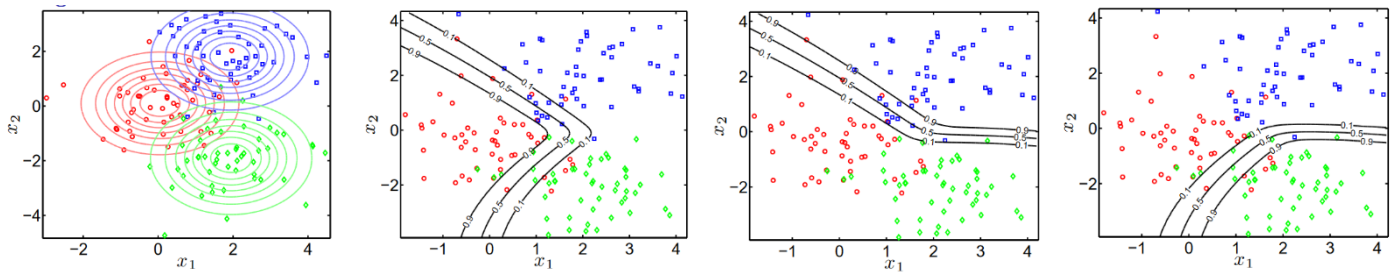
Naïve-Bayes

We add the assumption that all features of x are independent.

$$p(x) = p(x_d | x_{d-1}, \dots, x_1)p(x_{d-1} | x_{d-2}, \dots, x_1) \dots p(x_1) = p(x_d)p(x_{d-1}) \dots p(x_1)$$

Step1: fitting the class-conditional densities, i.e. find features' means and variances for each class.

Step2: evaluate densities at test point.



Performance Evaluation

0/1 loss

proportion of times classifier is wrong. Mean loss is defined as $\frac{1}{N} \sum \delta(t_n \neq t_n^*)$

- + Used for binary or multiclass classification.
- + Simple to compute and give a single value.
- Does not consider class imbalance

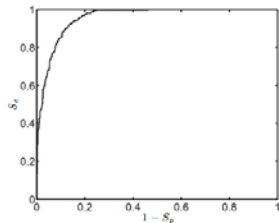
Sensitivity and specificity

True/False Negatives/Positives: true means correctly classified, and positives refer to class 1.

Sensitivity $S_e = \frac{TP}{TP+FN}$ (proportion of diseased classified as diseased).

Specificity $S_p = \frac{TN}{TN+FP}$ (proportion of healthy classified as healthy).

We would like both to be as high as possible



ROC Analysis

The Receiver Operating Characteristic curve shows how S_e and $1 - S_p$ vary as the threshold changes.

We can quantify the performance by computing the area under the ROC curve (AUC).

Confusion Matrix

		True Class	
		1	0
Predicted Class	1	TP	FP
	0	FN	TN

Lecture 5: Variational Inference

The aim of this part is to approximate the posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{t})$

1- Introduce $q(\mathbf{w})$

For simplicity $q(\mathbf{w}) = \prod q(w_i) = \prod \mathcal{N}(w_i|\mu_i, \sigma_i^2)$ (it's a choice and other alternatives can be used).

2- Distance between $q(\mathbf{w})$ and $p(\mathbf{w}|\mathbf{X}, \mathbf{t})$

A possible distance can be $KL[q(\mathbf{w})||p(\mathbf{w}|\mathbf{X}, \mathbf{t})] = \mathbb{E}_{q(\mathbf{w})} \left[\log \frac{q(\mathbf{w})}{p(\mathbf{w}|\mathbf{X}, \mathbf{t})} \right] = \mathbb{E}_{q(\mathbf{w})} [\log q(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w})} [\log p(\mathbf{w}|\mathbf{X}, \mathbf{t})]$

This is not a distance! Not symmetric! Does not satisfy triangular equality.

The second term is problematic because the posterior is intractable.

After rearranging: $\log p(\mathbf{t}|\mathbf{X}) - KL[q(\mathbf{w})||p(\mathbf{w}|\mathbf{X}, \mathbf{t})] = \mathbb{E}_{q(\mathbf{w})}[\log p(\mathbf{t}|\mathbf{w}, \mathbf{X})] - KL[q(\mathbf{w})||p(\mathbf{w})]$

We maximize the objective = $\mathbb{E}_{q(\mathbf{w})}[\log p(\mathbf{t}|\mathbf{w}, \mathbf{X})] - KL[q(\mathbf{w})||p(\mathbf{w})]$ wrt $q(\mathbf{w})$, i.e. change q params.

3- Optimization

If $p(\mathbf{w}) = \prod \mathcal{N}(w_i|0, s^2)$ we obtain $KL[q(\mathbf{w})||p(\mathbf{w})] = \frac{1}{2} \sum \left[\log \frac{s^2}{\sigma_i^2} - 1 + \frac{\sigma_i^2}{s^2} + \frac{\mu_i^2}{s^2} \right]$

$\mathbb{E}_{q(\mathbf{w})}[\log p(\mathbf{t}|\mathbf{w}, \mathbf{X})] = \int q(\mathbf{w}) \log p(\mathbf{t}|\mathbf{w}, \mathbf{X}) d\mathbf{w} \approx \frac{1}{N_{MC}} \sum_{h=1}^{N_{MC}} \log p(\mathbf{t}|\tilde{\mathbf{w}}^{(h)}, \mathbf{X})$ where $\tilde{\mathbf{w}}^{(h)}$ are sampled from $q(\mathbf{w})$ such that $(\tilde{\mathbf{w}}^{(h)})_i = \mu_i + \epsilon_i \sigma_i$.

Now we maximize $\widetilde{\text{objective}} = \frac{1}{N_{MC}} \sum_{h=1}^{N_{MC}} \log p(\mathbf{t}|\tilde{\mathbf{w}}^{(h)}, \mathbf{X}) - KL[q(\mathbf{w})||p(\mathbf{w})]$ with gradient-based optimization.
 $\text{vapr}' = \text{vpar} + \frac{\alpha_t}{2} \nabla_{\text{vpar}}(\widetilde{\text{objective}}); \alpha_t \rightarrow 0$

Lecture 6: Bayesian Unsupervised Learning

K-means

Each cluster is defined by a position in the input space $\boldsymbol{\mu}_k = [\mu_{k1}, \mu_{k2}]^T$.

Each \mathbf{x}_n is assigned to its closest cluster.

Euclidean distance is usually used: $d_{nk} = (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k)$

There is no analytical solution, so we use an iterative algorithm:

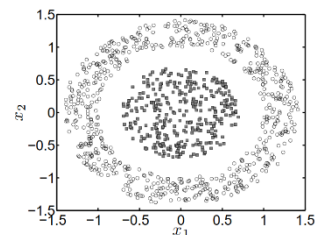
1. Randomly pick $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$
2. Assign each \mathbf{x}_n to the closest $\boldsymbol{\mu}_k$
3. Define $z_{nk} = 1$ if \mathbf{x}_n is assigned to $\boldsymbol{\mu}_k$
4. Update $\boldsymbol{\mu}_k$ to the average of \mathbf{x}_n 's $\boldsymbol{\mu}_k = \frac{\sum z_{nk} \mathbf{x}_n}{\sum z_{nk}}$
5. Return to 2 until assignment do not change

Kernelizing K-means

$$\begin{aligned} d_{nk} &= (\mathbf{x}_n - \boldsymbol{\mu}_k)^T (\mathbf{x}_n - \boldsymbol{\mu}_k) = \left(\mathbf{x}_n - \frac{\sum z_{nk} \mathbf{x}_n}{\sum z_{nk}} \right)^T \left(\mathbf{x}_n - \frac{\sum z_{nk} \mathbf{x}_n}{\sum z_{nk}} \right) \\ &= \left(\mathbf{x}_n - N_k^{-1} \sum z_{nk} \mathbf{x}_n \right)^T \left(\mathbf{x}_n - N_k^{-1} \sum z_{nk} \mathbf{x}_n \right) \\ &= \mathbf{x}_n^T \mathbf{x}_n - 2N_k^{-1} \sum_m z_{mk} \mathbf{x}_m^T \mathbf{x}_n + N_k^{-2} \sum_{m,l} z_{mk} z_{lk} \mathbf{x}_m^T \mathbf{x}_l \\ &= k(\mathbf{x}_n, \mathbf{x}_n) - 2N_k^{-1} \sum_m z_{mk} k(\mathbf{x}_n, \mathbf{x}_m) + N_k^{-2} \sum_{m,l} z_{mk} z_{lk} k(\mathbf{x}_m, \mathbf{x}_l) \end{aligned}$$

Algorithm:

1. Choose a kernel and any necessary params
2. Start with random assignments z_{nk}
3. For each \mathbf{x}_n assign it to the nearest center using the new distance
4. Return to 3 until assignment do not change



Mixture Models

Can we hypothesize a function that could have generated the data?