

Tutorial 10 Solutions

Sunday 13th November, 2016

1 Gaussian Discriminant Analysis

1.1 Quadratic separating surface

- Consider the following Gaussian Discriminant Classifier discussed in class. Let us say we have K classes with a multivariate Gaussian Model $\mathcal{N}(\mu_i, \Sigma_i)$ fitted for each class:

$$P(\phi(x)|C_1) = \mathcal{N}(\mu_1, \Sigma_1)$$

$$P(\phi(x)|C_i) = \mathcal{N}(\mu_i, \Sigma_i)$$

$$P(\phi(x)|C_K) = \mathcal{N}(\mu_K, \Sigma_i)$$

- Assumption: $\phi(x)$ is generated using **exactly one** $\mathcal{N}(\mu_i, \Sigma_i)$
- In the case of $K = 2$, with $P(C_1)$ and $P(C_2)$ known, separating surface will be $\{\phi(x) \mid P(C_1|\phi(x)) = P(C_2|\phi(x))\}$.

Prove that the separating surface will be **quadratic**.

ANSWER:

- If $\phi(x) \sim \mathcal{N}(\mu_i, \Sigma_i)$ (where $\phi(x) \in \mathfrak{R}^m$) then

$$p(\phi(x) | C_i) = \frac{1}{(2\pi)^{\frac{m}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp\left\{-\frac{(\phi(x) - \mu_i)^T \Sigma_i^{-1} (\phi(x) - \mu_i)}{2}\right\}$$

- So, the separating surface is $\phi(x)$ such that $\{\phi(x) \mid P(C_1|\phi(x)) = P(C_2|\phi(x))\} \iff \{\phi(x) \mid P(\phi(x) | C_1)P(C_1) = P(\phi(x) | C_2)P(C_2)\} \iff$ after taking logs, $\phi(x)$ such that

$$-(\phi(x) - \mu_1)^T \Sigma_1^{-1} (\phi(x) - \mu_1) + (\phi(x) - \mu_2)^T \Sigma_2^{-1} (\phi(x) - \mu_2) = b$$

where b contains terms independent of $\phi(x)$.

- This is indeed a **QUADRATIC equation!**

1.2 Linear Discriminant Analysis

- Now consider the following variant of the above Gaussian Discriminant Classifier. Let us say we have K classes with a multivariate Gaussian Model $\mathcal{N}(\mu_i, \Sigma)$ fitted for each class. That is, the covariance matrix Σ is now shared across the classes:

$$P(\phi(x)|C_1) = \mathcal{N}(\mu_1, \Sigma)$$

$$P(\phi(x)|C_i) = \mathcal{N}(\mu_i, \Sigma)$$

$$P(\phi(x)|C_K) = \mathcal{N}(\mu_K, \Sigma)$$

- Assumption: $\phi(x)$ is generated using **exactly one** $\mathcal{N}(\mu_i, \Sigma)$.
- As before, in the case of $K = 2$, with $P(C_1)$ and $P(C_2)$ known, separating surface will be $\{\phi(x) \mid P(C_1|\phi(x)) = P(C_2|\phi(x))\}$.

1. Q: Prove that the separating surface will now be **linear**.

Answer:

- If $\phi(x) \sim \mathcal{N}(\mu_i, \Sigma)$ (where $\phi(x) \in \mathfrak{R}^m$) then

$$p(\phi(x) \mid C_i) = \frac{1}{(2\pi)^{\frac{m}{2}} |\Sigma|^{\frac{1}{2}}} \exp \frac{-(\phi(x) - \mu_i)^T \Sigma^{-1} (\phi(x) - \mu_i)}{2}$$

- So, the separating surface is $\phi(x)$ such that $\{ \phi(x) \mid P(C_1|\phi(x)) = P(C_2|\phi(x)) \}$
 $\Leftrightarrow \{ \phi(x) \mid P(\phi(x) \mid C_1)P(C_1) = P(\phi(x) \mid C_2)P(C_2) \} \Leftrightarrow$ after taking logs,
 $\phi(x)$ such that

$$-(\phi(x) - \mu_1)^T \Sigma^{-1} (\phi(x) - \mu_1) + (\phi(x) - \mu_2)^T \Sigma^{-1} (\phi(x) - \mu_2) = b$$

where b contains terms independent of $\phi(x)$

\Leftrightarrow ...the above expression can actually be simplified by canceling out the terms involving $\phi^T(x)\Sigma^{-1}\phi(x)$

$$-\phi^T(x)\Sigma^{-1}\phi(x) \dots + \phi^T(x)\Sigma^{-1}\phi(x) = b$$

to finally give

$$2\phi^T(x)\Sigma^{-1}\mu_1 + 2\phi^T(x)\Sigma^{-1}\mu_2 - \mu_1^T(x)\Sigma^{-1}\mu_1 - \mu_2^T(x)\Sigma^{-1}\mu_2 = b$$

that is,

$$2\phi^T(x)\Sigma^{-1}\mu_1 + 2\phi^T(x)\Sigma^{-1}\mu_2 - \mu_1^T(x)\Sigma^{-1}\mu_1 = b'$$

which is a **LINEAR equation!** Here, $b' = \mu_1^T(x)\Sigma^{-1}\mu_1 + \mu_2^T(x)\Sigma^{-1}\mu_2 + b$ is independent of $\phi(x)$.

2. Q: What will be the maximum likelihood estimates for μ_i and Σ in this new case of different means but shared covariance matrix?

ANSWER: The Maximum Likelihood estimate for $\hat{\mu}_i$ will be the same as that for the Quadratic Discriminant Analysis, but that for a shared and single covariance estimate $\hat{\Sigma}$ will correspond to the average of covariance matrix estimates across examples from all the classes

$$\hat{\mu}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \phi(x_j^i) \rightarrow \text{Same as for QDA}$$

Earlier for QDA
 $\hat{\Sigma}_i = \frac{1}{n_i} \sum_{x \in C_i} (\phi(x) - \mu_i)(\phi(x) - \mu_i)^T$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^K \sum_{j=1}^{n_i} (\phi(x_j^i) - \mu_i)(\phi(x_j^i) - \mu_i)^T \rightarrow \text{X}$$

2 EM Algorithm for Mixture of Gaussians

Q: Show that the following algorithm for estimating the mean μ_i , the covariance matrix Σ_i and mixture components π_i for a mixture of Gaussians is an instance of the general EM algorithm

ANSWER: Has been discussed in Lecture 27.

Initialize $\mu_i^{(0)}$ to different random values and $\Sigma_i^{(0)}$ to I . Now iterate between the following

E Step and M Steps:

E Step:

1. For the posterior $p(z_i | \phi(x_j), \mu, \Sigma)$

$$p^{(t+1)}(z_i | \phi(x_j), \theta) = \frac{\pi_i \mathcal{N}(\phi(x); \mu_i^{(t)}, \Sigma_i^{(t)})}{\sum_{l=1}^K \pi_l \mathcal{N}(\phi(x); \mu_l^{(t)}, \Sigma_l^{(t)})}$$

M Steps:

1. For the prior π_i

$$\pi_i^{(t+1)} = \frac{1}{n} \sum_{j=1}^n p^{(t+1)}(z_i | \phi(x_j), \theta)$$

2. For μ_i

$$\mu_i^{(t+1)} = \frac{\sum_{j=1}^n p^{(t+1)}(z_i | \phi(x_j), \theta) \phi(x_j)}{\sum_{j=1}^n p^{(t+1)}(z_i | \phi(x_j), \theta)}$$

3. For Σ_i

$$\underline{\Sigma}_i^{(t+1)} = \frac{\sum_{j=1}^n p^{(t+1)}(z_i | \phi(x_j), \theta) (\phi(x_j) - \mu_i^{(t+1)}) (\phi(x_j) - \mu_i^{(t+1)})^T}{\sum_{j=1}^n p^{(t+1)}(z_i | \phi(x_j), \theta)}$$

Q: Note that this algorithm is for the Mixture of Gaussians assuming a different ~~Σ~~ covariance matrix Σ_i for each class C_i . What will be the algorithm like, if we assume a shared covariance matrix Σ across all classes (that is, the Linear Discriminant Analysis discussed in Section 1.2)?

ANSWER: We will simply build on the solution to the Linear Discriminant case from Section 2.1 and simply replace multiple class-specific estimates Σ_i with a single estimate Σ :

Mstep: $\rightarrow \Sigma^{(t+1)} = \frac{1}{m} \sum_x \sum_i q_i^t(z_i|x) (\phi(x) - \mu_i^{t+1}) (\phi(x) - \mu_i^{t+1}) / \sum_x \sum_i q_i^t(z_i|x)$ (based on (*))

3 Convergence of Hard K-Means Algorithm

Prove the following claim: The K-Means Clustering algorithm will converge in a finite number of iterations.

1. **Proof Sketch:** At each iteration, the K-Means algorithm reduces the objective $\sum_{j=1}^m \sum_{l=1}^K P_{l,j} \|\phi(\mathbf{x}^{(j)}) - \mu_l\|^2$ and stops when this objective does not reduce any further.

2. Hint1: $P^{(t+1)} = \operatorname{argmin}_P \sum_{j=1}^m \sum_{l=1}^K P_{l,j} \|\phi(\mathbf{x}^{(j)}) - \mu_l^{(t)}\|^2$

(E step equivalent)

3. Hint2: $\mu^{(t+1)} = \operatorname{argmin}_\mu \sum_{j=1}^m \sum_{l=1}^K P_{l,j}^{(t+1)} \|\phi(\mathbf{x}^{(j)}) - \mu_l\|^2$

(M step equivalent)

4. Hint3: Only a finite number of combinations of $P_{l,j}$ are possible.

4 (Optional) Bayesian Inference from Multinomial to Naive Bayes (Useful for text data)

First we summarize the conjugate prior, MLE and Bayesian estimation for Multinomial

- In the case of the Multinomial distribution (extension to the binomial distribution) a variable X could assume one of t possible values $V_1, V_2 \dots V_t$ with parameters $\Pr(X = V_j) = \mu_j$. Each observation \mathbf{X}_k (for $k \in [1, n]$) is modeled as a vector $\mathbf{X}_k = [X_{k,1} \dots X_{k,j} \dots X_{k,t}]$ with $X_{k,j} = 1$ if and only if value of \mathbf{X}_k was observed to be V_j and $X_{k,j} = 0$ otherwise. Eg: In the case of the toss of dice, $t = 6$.

- The maximum likelihood estimate for the mean is given by:

$$\hat{\mu}_j = \frac{\sum_{k=1}^n X_{k,j}}{n} = \frac{n_j}{n} \quad (1)$$

where, given n iid observations of a multinomial random variable X , $n_j = \sum_{k=1}^n X_{k,j}$ is the number of times $X = V_j$ was observed,

- Conjugate prior follows $\text{Dir}(\alpha_1 \dots \alpha_n)$

- Posterior is $\text{Dir}(\dots \alpha_l + \sum_{k=1}^n X_{k,j} \dots)$

- The expectation of μ under $\text{Dir}(\alpha_1 \dots \alpha_n)$ is given by:

$$E[\mu]_{\text{Dir}(\alpha_1 \dots \alpha_n)} = \left[\frac{\alpha_1}{\sum \alpha_l} \dots \frac{\alpha_l}{\sum \alpha_l} \right] \quad (2)$$

- The (posterior) expectation of μ under $\text{Dir}(\dots \alpha_j + \sum_{k=1}^n X_{k,j} \dots)$ is given by:

$$E[\mu]_{\text{Dir}(\dots \alpha_j + \sum_{k=1}^n X_{k,j} \dots)} = \left[\frac{\alpha_1 + \sum_k X_{k,1}}{\sum \alpha_l + n} \dots \frac{\alpha_j + \sum_k X_{k,j}}{\sum \alpha_l + n} \dots \right] \quad (3)$$

Now recall that we extended single Multinomial to Multinomial Naive Bayes that has class conditioned independent features, each of which is Multinomial. We also discussed Maximum Likelihood estimation of Multinomial Naive Bayes. All of that is summarized below:

- $\langle X_k, C_i \rangle$: Tuple with example X_k belonging to class C_i . $\text{Pr}(C_i)$ is prior probability of class C_i .
- $\phi_1(X_k), \dots, \phi_m(x_k)$: The feature vector for X_k
- $P(\phi_q(x)|C_i) \sim \text{Mult}(\mu_{1,i}^q \dots \mu_{t_q,i}^q)$; that is, each feature ϕ_q follows multinomial distribution Bayes
 1. $[V_1^1 \dots V_{t_1}^1] \dots [V_1^q \dots V_{t_q}^q] \dots [V_1^m \dots V_{t_m}^m]$: Set of values that could be taken by each of $\phi_1, \phi_2 \dots \phi_m$ respectively
 2. $[\mu_{1,i}^1 \dots \mu_{t_1,i}^1] \dots [\mu_{1,i}^q \dots \mu_{t_q,i}^q] \dots [\mu_{1,i}^m \dots \mu_{t_m,i}^m]$: Parameters for each of $\phi_1, \phi_2 \dots \phi_m$ respectively for class C_i
- $P(\phi_1(x) \dots \phi_m(x)|C_i) = \prod_{q=1}^m P(\phi_q(x)|C_i)$: Feature are independent given the class
- Maximum Likelihood Estimate for Naive Bayes: Let

$\#C_i$ = No. of times $c(X_k) = C_i$ across all k 's in the dataset

$n_{j,i}^q$ = No. of times $\phi_q(X_k) = V_j$ and $c(X_k) = C_i$ across all the k 's

$$n_{j,i}^q = \sum_k \delta(\phi_q(X_k), V_j^q) \delta(c(X_k), C_i)$$

then

$$\hat{\mu}_{j,i}^q = \frac{n_{j,i}^q}{\sum_{j'=1}^n n_{j',i}^q}$$

$$\widehat{Pr}_{c_i} = \frac{\#C_i}{\sum_{i'} \#C_{i'}}$$

Q: Extend Bayesian Inference using the Dirichlet prior from Multinomial to Naive Bayes

ANSWER: The setting for Naive Bayes with Dirichlet prior on Multivariate Bernoulli distribution is as follows

- For each data point X_k which belongs to class C_i there are a set of m features given by $\phi_1(X_k) \dots \phi_q(X_k) \dots \phi_m(X_k) | C_i$
- Each parameter (vector) $\mu_i^q = [\mu_{1,i}^q \dots \mu_{j,i}^q \dots \mu_{t_q,i}^q]$ corresponding to a feature $\phi_q(\cdot)$ has a probability distribution given by

$$p(\mu_i^1) \sim Dir(a_{1,i}^1 \dots a_{j,i}^1 \dots a_{t_1,i}^1) \dots$$

$$p(\mu_i^q) \sim Dir(a_{1,i}^q \dots a_{j,i}^q \dots a_{t_q,i}^q) \dots$$

$$p(\mu_i^m) \sim Dir(a_{1,i}^m \dots a_{j,i}^m \dots a_{t_m,i}^m)$$
 where $a_{j,i}^q$ are the multivariate dirichlet prior parameters
- Let $n_{j,i}^q$ be the number of times attribute $\phi_q(\cdot)$ was observed with value V_j^q amongst examples belonging to class C_i .
- Then, from our previous analysis of Bayesian estimation for Multinomial¹

$$p(\mu_i^1 | D) \sim Dir(a_{1,i}^1 + n_{1,i}^1 \dots a_{j,i}^1 + n_{j,i}^1 \dots a_{t_1,i}^1 + n_{t_1,i}^1) \dots$$

$$p(\mu_i^q | D) \sim Dir(a_{1,i}^q + n_{1,i}^q \dots a_{j,i}^q + n_{j,i}^q \dots a_{t_q,i}^q + n_{t_q,i}^q) \dots$$

$$p(\mu_i^m | D) \sim Dir(a_{1,i}^m + n_{1,i}^m \dots a_{j,i}^m + n_{j,i}^m \dots a_{t_m,i}^m + n_{t_m,i}^m).$$

¹See page 12 of <http://23.253.82.180/course/307/865/2238>.