

QUESTION

1

Given a soft margin linear SVM problem

$$\max_{\mathbf{0} \leq \boldsymbol{\alpha} \leq C} f(\boldsymbol{\alpha}) = \boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{G} \boldsymbol{\alpha}$$

where all the vectors and matrices have their usual meanings as discussed in class. The given problem handles optimization process using a variant of coordinate gradient descent based approach. The updation rule for this approach is as follows:

$$\begin{aligned}\alpha_n &= \alpha_n + \delta_* \\ \delta_* &= \arg \max_{\delta} f(\boldsymbol{\alpha} + \delta \mathbf{e}_n)\end{aligned}$$

where \mathbf{e}_n denotes a vector of all zeros except a 1 entry at n.

$$\begin{aligned}f(\boldsymbol{\alpha} + \delta \mathbf{e}_n) &= (\boldsymbol{\alpha} + \delta \mathbf{e}_n)^T \mathbf{1} - \frac{1}{2} (\boldsymbol{\alpha} + \delta \mathbf{e}_n)^T \mathbf{G} (\boldsymbol{\alpha} + \delta \mathbf{e}_n) \\ &= f(\boldsymbol{\alpha}) + \delta (1 - \mathbf{e}_n^T \mathbf{G} \boldsymbol{\alpha}) - \frac{1}{2} \delta^2 \mathbf{e}_n^T \mathbf{G} \mathbf{e}_n\end{aligned}$$

Maximizing w.r.t. to δ we have:

$$\begin{aligned}\frac{\partial}{\partial \delta} f(\boldsymbol{\alpha} + \delta \mathbf{e}_n) &= 1 - \mathbf{e}_n^T \mathbf{G} \boldsymbol{\alpha} - \delta \mathbf{e}_n^T \mathbf{G} \mathbf{e}_n = 0 \\ \delta^* &= \frac{(1 - \mathbf{e}_n^T \mathbf{G} \boldsymbol{\alpha})}{\mathbf{e}_n^T \mathbf{G} \mathbf{e}_n}\end{aligned}$$

Simplifying further we have

$$\delta^* = \frac{1 - y_n \mathbf{w}^T \mathbf{x}_n}{\|\mathbf{x}_n\|^2}$$

$$\because \mathbf{e}_n^T \mathbf{G} \boldsymbol{\alpha} = \sum_{m=1}^N y_m \mathbf{x}_n^T \mathbf{x}_m \alpha_m = y_n x_n^T \left(\sum_{m=1}^N y_m \mathbf{x}_m \alpha_m \right) = y_n \mathbf{w}^T \mathbf{x}_n \text{ and } \mathbf{e}_n^T \mathbf{G} \mathbf{e}_n = y_n^2 \|\mathbf{x}_n\|^2 = \|\mathbf{x}_n\|^2$$

The above method doesn't put a constraint on updated α_n . To handle this we use projected gradient descent i.e if $\alpha_n + \delta^* \in [0, C]$ we update in a regular fashion else if $\alpha_n + \delta^* < 0$ we update α_n to 0 else if $\alpha_n + \delta^* > C$ we update α_n to C .

Algorithm 1 Co-ordinate Ascent for Soft Margin SVM

1. Initialize $\boldsymbol{\alpha}$ as $\hat{\boldsymbol{\alpha}}$ s.t $0 < \alpha_n < C \forall n = 1, \dots, N$
2. For i in $[1, \dots, N]$

$$\alpha_i = \begin{cases} \alpha_i + \delta^* & \alpha_i + \delta^* \in [0, C] \\ 0 & \alpha_i + \delta^* < 0 \\ C & \alpha_i + \delta^* > C \end{cases}$$

3. If not converged go to step 2

Ridge Regression: Revisited Question 2

- Recall the ridge regression problem

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \mathbf{w}^\top \mathbf{w}$$

- The solution to this problem was

$$\mathbf{w} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D \right) \left(\sum_{n=1}^N y_n \mathbf{x}_n \right) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Inputs don't appear as inner-products here . They actually do! :-)
- Matrix inversion lemma: $(\mathbf{F}\mathbf{H}^{-1}\mathbf{G} - \mathbf{E})^{-1}\mathbf{F}\mathbf{H}^{-1} = \mathbf{E}^{-1}\mathbf{F}(\mathbf{G}\mathbf{E}^{-1}\mathbf{F} - \mathbf{H})^{-1}$
- The lemma allows us to rewrite \mathbf{w} as

$$\mathbf{w} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = \mathbf{X}^\top \boldsymbol{\alpha} = \sum_{n=1}^N \alpha_n \mathbf{x}_n$$

where $\boldsymbol{\alpha} = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$ is an $N \times 1$ vector of **dual variables**, and $K_{nm} = \mathbf{x}_n^\top \mathbf{x}_m$

- Note: $\mathbf{w} = \sum_{n=1}^N \alpha_n \mathbf{x}_n$ is known as "dual" form of ridge regression solution. However, so far it is still a linear model. But now it is easily kernelizable.

Gaussian Distribution Question 3

- The (multivariate) Gaussian with mean μ and cov. matrix Σ

$$\begin{aligned}\mathcal{N}(x|\mu, \Sigma) &= \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp \left\{ -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right\} \\ &= \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp \left\{ -\frac{1}{2} \text{trace} [\Sigma^{-1} S] \right\} \quad \text{where } S = (x - \mu)(x - \mu)^\top\end{aligned}$$

- An alternate representation: The “information form”

$$\mathcal{N}_c(x|\xi, \Lambda) = (2\pi)^{-D/2} |\Lambda|^{1/2} \exp \left\{ -\frac{1}{2} (x^\top \Lambda x + \xi^\top \Lambda^{-1} \xi - 2x^\top \xi) \right\}$$

where $\Lambda = \Sigma^{-1}$ and $\xi = \Sigma^{-1}\mu$ are the “natural parameters” (recall exp. family).

- Note that there is a term quadratic in x (involves $\Lambda = \Sigma^{-1}$) and linear in x (involves $\xi = \Sigma^{-1}\mu$)
- Information form can help recognize μ and Σ of a Gaussian when doing algebraic manipulations

Estimating Parameters of Gaussian: MLE

- Given: N i.i.d. observations $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ from a multivariate Gaussian $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- Goal: Estimate $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$. Simple to do MLE for this task

$$\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \log p(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Plugging in $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\}$ and ignoring the constants

$$\begin{aligned}\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{N}{2} \log |\boldsymbol{\Sigma}|^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \\ &= \frac{N}{2} \log |\boldsymbol{\Sigma}|^{-1} - \frac{1}{2} \sum_{n=1}^N \text{trace}[\boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top] \\ &= \frac{N}{2} \log |\boldsymbol{\Sigma}|^{-1} - \frac{1}{2} \text{trace}[\boldsymbol{\Sigma}^{-1} \mathbf{S}_\mu] \quad \left[\text{where } \mathbf{S}_\mu = \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top \right]\end{aligned}$$

Estimating Parameters of Gaussian: MLE

- Taking (partial) derivatives w.r.t. μ and setting to zero

$$\frac{\partial}{\partial \mu} \mathcal{L}(\mu, \Sigma) = \frac{\partial}{\partial \mu} \left[\frac{N}{2} \log |\Sigma|^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)^\top \Sigma^{-1} (\mathbf{x}_n - \mu) \right] = -\frac{1}{2} \sum_{n=1}^N (\Sigma^{-1} + \Sigma^{-\top}) (\mathbf{x}_n - \mu) = 0$$

which gives the following MLE solution for the multivariate Gaussian's mean

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- Taking derivatives w.r.t. $\Lambda = \Sigma^{-1}$ (instead of Σ ; [leads to simpler derivatives](#)) and setting to zero

$$\frac{\partial}{\partial \Lambda} \mathcal{L}(\mu, \Lambda) = \frac{\partial}{\partial \Lambda} \left[\frac{N}{2} \log |\Lambda| - \frac{1}{2} \text{trace}[\Lambda \mathbf{S}_\mu] \right] = \frac{N}{2} \Lambda^{-\top} - \frac{1}{2} \mathbf{S}_\mu^\top = \frac{N}{2} \Lambda^{-1} - \frac{1}{2} \mathbf{S}_\mu = \frac{N}{2} \Sigma - \frac{1}{2} \mathbf{S}_\mu = 0$$

which gives the following MLE solution for the multivariate Gaussian's covariance matrix

$$\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{ML})(\mathbf{x}_n - \mu_{ML})^\top$$



Question 4

$$y = \{y_1, y_2, \dots, y_n\}$$
$$P(y|\lambda) = \frac{\lambda^{y_n} e^{-\lambda}}{y!} \quad (\text{Poisson})$$

$$P(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \quad (\text{Gamma})$$

① MLE and MAP:

$$\hat{\lambda}_{MLE} = \underset{\lambda}{\operatorname{argmax}} \sum_{n=1}^N \log P(y_n|\lambda) \quad (\text{i.i.d. data})$$

$$\sum_{n=1}^N \log P(y_n|\lambda) = \left[\sum_{n=1}^N y_n \log \lambda \right] - N \cancel{\lambda} + \text{(constant)} \quad (\text{terms that don't depend on } \lambda)$$

Taking derivative and setting it to zero

$$\sum y_n \times \frac{1}{\lambda} - N = 0$$

$$\Rightarrow \boxed{\hat{\lambda}_{MLE} = \frac{\sum_{n=1}^N y_n}{N}}$$

MAP estimation will be almost identical with the extra $\log P(\lambda)$ term.

$$\hat{\lambda}_{MAP} = \underset{\lambda}{\operatorname{argmax}} \left[\sum_{n=1}^N \log P(y_n|\lambda) + \log P(\lambda) \right]$$

$$\log P(\lambda) = (\alpha-1) \log \lambda - \beta \lambda + \text{Constant}$$

terms that
don't depend on λ

The MAP objective will be

$$\sum_{n=1}^N y_n \log \lambda - N\lambda + (\alpha-1) \log \lambda - \beta \lambda$$

Maximizing w.r.t. λ will give the MAP solution

$$\hat{\lambda}_{\text{MAP}} = \frac{\sum_{n=1}^N y_n + \alpha-1}{N+\beta}$$

② Posterior distribution of λ

$$P(\lambda|y) = \frac{P(\lambda) P(y|\lambda)}{P(y)} = P(\lambda) \prod_{n=1}^N P(y_n|\lambda)$$

Since the prior (Gamma) and the likelihood (Poisson) ~~are~~ are conjugate, the posterior must be gamma!

Let's multiply the terms $P(\lambda)$ and $P(y_n|\lambda)$
are try to "identify" this gamma distribution's parameters.

$$P(\lambda|y) \propto P(\lambda) \prod_{n=1}^N P(y_n|\lambda)$$

$$P(\lambda|y) \propto \lambda^{d-1} e^{-\beta\lambda} \times \prod_{n=1}^N \frac{y_n}{\lambda} e^{-\lambda}$$

$$\propto \lambda^{\sum_{n=1}^N y_n + d-1} e^{-(\beta+N)\lambda}$$

(ignoring the terms that don't depend on λ)

The above expression is clearly in form of a gamma distribution with

$$\text{Shape} = \sum_{n=1}^N y_n + d-1 \quad \left[\begin{array}{l} \text{No need to worry} \\ \text{about constant} \\ \text{of proportionality.} \\ \text{It must be a gamma} \end{array} \right]$$

$$\text{Rate} = \beta + N$$

Thus $P(y|\lambda) = \text{Gamma}\left(\sum_{n=1}^N y_n + d-1, \beta + N\right)$

$$(3) \quad \hat{\lambda}_{MAP} = \frac{\sum_{n=1}^N y_n + d-1}{N+\beta} = \frac{\sum_{n=1}^N y_n}{N+\beta} + \frac{(d-1)}{(N+\beta)}$$

Sum to 1
(convex comb.)

$$= \left[\frac{\sum_{n=1}^N y_n}{N} \right] \times \frac{N}{N+\beta} + [d-1] \times \frac{\beta}{N+\beta}$$

\downarrow MLE \downarrow Prior's mode

Posterior's mean:

$$\frac{\sum_{n=1}^N y_n + d-1}{N+\beta} = \frac{\sum_{n=1}^N y_n}{N} \times \frac{N}{N+\beta} + \left[\frac{d-1}{\beta} \right] \times \frac{\beta}{N+\beta}$$

\downarrow MLE \downarrow Prior's mean

$$\begin{aligned}
 ④ P(y_*|y) &= \int P(y_*, \lambda | y) d\lambda \\
 &= \int P(y_*|\lambda) P(\lambda|y) d\lambda \\
 &\approx P(y_*|\hat{\lambda}_{MLE}) \quad (\text{if using MLE}) \\
 &\quad \text{or} \\
 &P(y_*|\hat{\lambda}_{MAP}) \quad (\text{if using MAP})
 \end{aligned}$$

In both these cases, $P(y_*|y)$ is simply a Poisson with parameters λ_{MLE} or λ_{MAP}

If using the full posterior, which is $P(\lambda|y) = \text{Gamma}(\sum_{n=1}^N y_n + \alpha, \beta + N)$, $P(y_*|y)$ will be

$$\begin{aligned}
 P(y_*|y) &= \int P(y_*|\lambda) P(\lambda|y) d\lambda \\
 &= \int \frac{\lambda^{y_*} e^{-\lambda}}{y_*!} \times \frac{(\beta + N)^{\sum_{n=1}^N y_n + \alpha}}{\prod_{n=1}^N y_n!} \lambda^{\sum_{n=1}^N y_n + \alpha - 1} e^{-(\beta + N)} d\lambda
 \end{aligned}$$

The above is actually a mixture of infinite many Poisson distributions. The result is actually not a Poisson but ~~is~~ the "Negative Binomial" distribution.

(if you are interested in knowing more about this result, you may refer to the wikipedia article of NB distribution).

(This part was just for your "General knowledge" :-))

Question 5

$$X = \{x_1, x_2, \dots, x_N\}$$

$$Y = \{y_1, y_2, \dots, y_N\}$$

The likelihood (both X and Y are modeled here)

$$P(X, Y | \Theta) = \prod_{n=1}^N P(x_n, y_n | \Theta)$$

$$\Theta = \{\mu_k, \Sigma_k, \pi_k\}_{k=1}^K = \prod_{n=1}^N \prod_{k=1}^K [P(x_n, y_n=k | \Theta)]$$

Note that only one term in the product will be selected based on the value of y_n

Thus log-likelihood will be

$$\log P(X, Y | \Theta) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] [\log P(x_n, y_n=k | \Theta)]$$

$$= \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] [\log P(x_n | y_n=k, \Theta) P(y_n=k | \Theta)]$$

$$L(\Theta) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] [\log N(x_n | \mu_k, \Sigma_k) + \log \pi_k]$$

To optimize w.r.t. π_k , we take partial derivatives w.r.t. π_k BUT need to use the constraint $\sum_k \pi_k = 1$. Also note that for $\{\pi_k\}_{k=1}^K$, the part of the objective that matters is

$$L(\pi_1, \dots, \pi_K) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] \log \pi_k$$

(Lagrangian form)
Need to minimize over π ,
max over y

$$\sum_{k=1}^K N_k \log \pi_k$$

\downarrow
 $\# \text{ of points}$
 $\text{with } y_n=k$

The Constrained formulation will be

$$L(\pi, \lambda) = \underbrace{\sum_{k=1}^K N_k \log \pi_k}_{\text{part of}} + \lambda \left[\sum_{k=1}^K \pi_k - 1 \right]$$

Taking derivative w.r.t. π_k ,

$$\frac{N_k}{\pi_k} + \lambda = 0$$

$$\pi_k = -\frac{N_k}{\lambda}$$

Also, using $\sum_{k=1}^K \pi_k = 1$ gives $\lambda = -N$

thus
$$\boxed{\pi_k = \frac{N_k}{N}}$$

Now, for μ_k and Σ_k

The relevant part of the objective function is

$$L(\{\mu_k, \Sigma_k\}) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}[y_n=k] \log \mathcal{N}(x_n | \mu_k, \Sigma_k)$$

MLE for μ_k, Σ_k is exactly the same as the MLE procedure for a multivariate Gaussian, but only using inputs x_n from class K .

(See attached slides for Gaussian MLE)

I have provided additional slides on the class webpage showing how to do it for multivariate Gaussian. The procedure for this problem will be exactly the same (but only using x_n with $y_n = k$)

Gaussian Distribution

- The (multivariate) Gaussian with mean μ and cov. matrix Σ

$$\begin{aligned}\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} \text{trace} [\boldsymbol{\Sigma}^{-1} \mathbf{S}] \right\} \quad \text{where } \mathbf{S} = (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top\end{aligned}$$

- An alternate representation: The “information form”

$$\mathcal{N}_c(\mathbf{x}|\boldsymbol{\xi}, \boldsymbol{\Lambda}) = (2\pi)^{-D/2} |\boldsymbol{\Lambda}|^{1/2} \exp \left\{ -\frac{1}{2} \left(\mathbf{x}^\top \boldsymbol{\Lambda} \mathbf{x} + \boldsymbol{\xi}^\top \boldsymbol{\Lambda}^{-1} \boldsymbol{\xi} - 2\mathbf{x}^\top \boldsymbol{\xi} \right) \right\}$$

where $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$ and $\boldsymbol{\xi} = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ are the “natural parameters” (recall exp. family).

- Note that there is a term quadratic in \mathbf{x} (involves $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$) and linear in \mathbf{x} (involves $\boldsymbol{\xi} = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$)
- Information form can help recognize $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ of a Gaussian when doing algebraic manipulations

Estimating Parameters of Gaussian: MLE

- Given: N i.i.d. observations $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ from a multivariate Gaussian $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- Goal: Estimate $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$. Simple to do MLE for this task

$$\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \log p(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Plugging in $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\}$ and ignoring the constants

$$\begin{aligned}\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{N}{2} \log |\boldsymbol{\Sigma}|^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \\ &= \frac{N}{2} \log |\boldsymbol{\Sigma}|^{-1} - \frac{1}{2} \sum_{n=1}^N \text{trace}[\boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top] \\ &= \frac{N}{2} \log |\boldsymbol{\Sigma}|^{-1} - \frac{1}{2} \text{trace}[\boldsymbol{\Sigma}^{-1} \mathbf{S}_\mu] \quad \left[\text{where } \mathbf{S}_\mu = \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top \right]\end{aligned}$$

Estimating Parameters of Gaussian: MLE

- Taking (partial) derivatives w.r.t. μ and setting to zero

$$\frac{\partial}{\partial \mu} \mathcal{L}(\mu, \Sigma) = \frac{\partial}{\partial \mu} \left[\frac{N}{2} \log |\Sigma|^{-1} - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)^\top \Sigma^{-1} (\mathbf{x}_n - \mu) \right] = -\frac{1}{2} \sum_{n=1}^N (\Sigma^{-1} + \Sigma^{-\top}) (\mathbf{x}_n - \mu) = 0$$

which gives the following MLE solution for the multivariate Gaussian's mean

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- Taking derivatives w.r.t. $\Lambda = \Sigma^{-1}$ (instead of Σ ; [leads to simpler derivatives](#)) and setting to zero

$$\frac{\partial}{\partial \Lambda} \mathcal{L}(\mu, \Lambda) = \frac{\partial}{\partial \Lambda} \left[\frac{N}{2} \log |\Lambda| - \frac{1}{2} \text{trace}[\Lambda \mathbf{S}_\mu] \right] = \frac{N}{2} \Lambda^{-\top} - \frac{1}{2} \mathbf{S}_\mu^\top = \frac{N}{2} \Lambda^{-1} - \frac{1}{2} \mathbf{S}_\mu = \frac{N}{2} \Sigma - \frac{1}{2} \mathbf{S}_\mu = 0$$

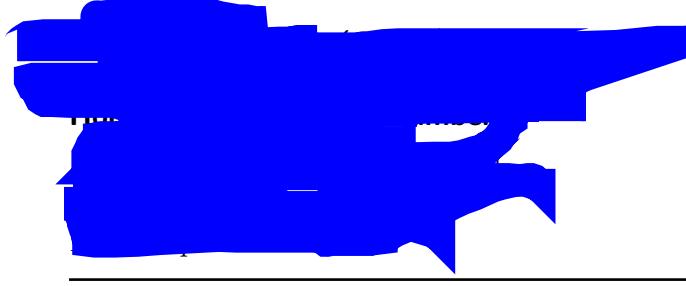
which gives the following MLE solution for the multivariate Gaussian's covariance matrix

$$\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{ML})(\mathbf{x}_n - \mu_{ML})^\top$$

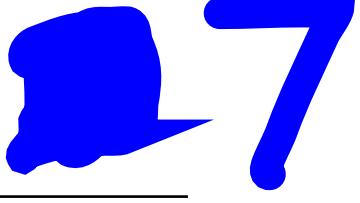


Question 6

Skipping the solution to practice problem 6 because we had already seen the Beta-Bernoulli case in class and the derivations for the Dirichlet-multinoulli case are almost identical.



QUESTION



Consider a logistic regression model $p(y_n|\mathbf{x}_n, \mathbf{w}) = \frac{1}{1+\exp(-y_n\mathbf{w}^T\mathbf{x}_n)}$, with a zero-mean Gaussian prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \lambda^{-1}\mathbf{I})$ and $y_n \in \{+1, -1\}$.

$$\therefore p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \frac{1}{1+\exp(-y_n\mathbf{w}^T\mathbf{x}_n)} \text{ (Assuming i.i.d.)}$$

For MAP estimation we have $\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{y}, \mathbf{X})$

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y}|\mathbf{w}, \mathbf{X})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

Since $p(\mathbf{y}|\mathbf{X})$ is independent of \mathbf{w} we have $\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} p(\mathbf{y}|\mathbf{w}, \mathbf{X})p(\mathbf{w})$

$$\begin{aligned} \hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} -\log(p(\mathbf{y}|\mathbf{w}, \mathbf{X})) - \log(p(\mathbf{w})) \\ &= \arg \min_{\mathbf{w}} -\log\left(\prod_{n=1}^N \frac{1}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)}\right) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \text{constant} \\ &= \arg \min_{\mathbf{w}} -\log\left(\prod_{n=1}^N \frac{1}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)}\right) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ &= \arg \min_{\mathbf{w}} -\sum_{n=1}^N \log\left(\frac{1}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)}\right) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

$$\text{Consider } \mathcal{L}(\mathbf{w}) = -\sum_{n=1}^N \log\left(\frac{1}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)}\right) + \frac{\lambda^2}{2} \mathbf{w}^T \mathbf{w}$$

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = -\sum_{n=1}^N \frac{\exp(-y_n \mathbf{w}^T \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)} y_n \mathbf{x}_n + \lambda \mathbf{w}$$

Setting $\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = 0$, we have

$$\begin{aligned} \hat{\mathbf{w}} &= \frac{1}{\lambda} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{w}^T \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)} y_n \mathbf{x}_n \\ &= \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \end{aligned}$$

$$\text{where } \alpha_n = \frac{1}{\lambda} \frac{\exp(-y_n \mathbf{w}^T \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)} = \frac{1}{\lambda} \frac{1}{1 + \exp(y_n \mathbf{w}^T \mathbf{x}_n)}.$$

α_n can also be written as $\frac{1}{\lambda}(1 - p(y_n|\mathbf{x}_n))$ which is nothing but the probability of getting incorrect label on input data \mathbf{x}_n scaled by a factor of $\frac{1}{\lambda}$. Now consider if a small probability of getting correct label with current \mathbf{w} the $\frac{1}{\lambda}(1 - p(y_n|\mathbf{x}_n))$ or α_n will be large contribution from this \mathbf{x}_n, y_n in calculating new $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ (taking a particular case of GD update when $\lambda = 1$). Now we attempt to show this change will improve \mathbf{w} w.r.t this particular \mathbf{x}_n . This is very similar to a perceptron update (mistake-driven), where the input having incorrect prediction will lead to significant change in weight vector. Also a factor of $\frac{1}{\lambda}$ will act as a regularizer to regulate the contribution of each \mathbf{x}_n .

QUESTION

8

For the given problem class-marginal distribution $p(y = 1) = \pi$ and each class-conditional distribution is defined as a product of D Bernoulli distributions, i.e., $p(\mathbf{x}|y = 1) = \prod_{d=1}^D p(x_d|y = 1)$

and $p(\mathbf{x}|y = 0) = \prod_{d=1}^D p(x_d|y = 0)$ where $p(x_d|y = 1) = Bernoulli(x_d|\mu_{d,1})$ and $p(x_d|y = 0) = Bernoulli(x_d|\mu_{d,0})$.

Assuming all the parameters have been estimated from the data, we can write the predictive distribution as follows

$$\begin{aligned} p(y = 1|\mathbf{x}) &= \frac{p(y = 1)p(\mathbf{x}|y = 1)}{p(y = 1)p(\mathbf{x}|y = 1) + p(y = 0)p(\mathbf{x}|y = 0)} \\ &= \frac{\pi \prod_{d=1}^D p(x_d|y = 1)}{\pi \prod_{d=1}^D p(x_d|y = 1) + (1 - \pi) \prod_{d=1}^D p(x_d|y = 0)} \\ &= \frac{1}{1 + \frac{(1 - \pi) \prod_{d=1}^D p(x_d|y = 0)}{\pi \prod_{d=1}^D p(x_d|y = 1)}} \\ &= \frac{1}{1 + \frac{1 - \pi \prod_{d=1}^D \mu_{d,0}^{x_d} (1 - \mu_{d,0})^{1-x_d}}{\pi \prod_{d=1}^D \mu_{d,1}^{x_d} (1 - \mu_{d,1})^{1-x_d}}} \end{aligned}$$

Comparing the above expression with $\sigma(\mathbf{w}^T \mathbf{x})$ i.e. $\exp(-(\mathbf{w}^T \mathbf{x} + b)) = \frac{1-\pi}{\pi} \frac{\prod_{d=1}^D \mu_{d,0}^{x_d} (1 - \mu_{d,0})^{1-x_d}}{\prod_{d=1}^D \mu_{d,1}^{x_d} (1 - \mu_{d,1})^{1-x_d}}$

we have

$$\begin{aligned} -\mathbf{w}^T \mathbf{x} - b &= \log\left(\frac{1 - \pi}{\pi}\right) + \sum_{d=1}^D \log(\mu_{d,0}^{x_d} (1 - \mu_{d,0})^{1-x_d}) - \sum_{d=1}^D \log(\mu_{d,1}^{x_d} (1 - \mu_{d,1})^{1-x_d}) \\ &= \log\left(\frac{1 - \pi}{\pi}\right) + \sum_{d=1}^D x_d \log \mu_{d,0} + \sum_{d=1}^D (1 - x_d) \log(1 - \mu_{d,0}) - \sum_{d=1}^D x_d \log \mu_{d,1} \\ &\quad - \sum_{d=1}^D (1 - x_d) \log(1 - \mu_{d,1}) \\ &= \sum_{d=1}^D x_d \log\left(\frac{\mu_{d,0}(1 - \mu_{d,1})}{(1 - \mu_{d,0})\mu_{d,1}}\right) + \log\left(\frac{1 - \pi}{\pi} \prod_{d=1}^D \frac{1 - \mu_{d,0}}{1 - \mu_{d,1}}\right) \end{aligned}$$

From above we can see that we can find an equivalent probabilistic discriminative model for this problem with

$$\mathbf{w} = \begin{bmatrix} \log\left(\frac{\mu_{1,1}(1-\mu_{1,0})}{(1-\mu_{1,1})\mu_{1,0}}\right) \\ \log\left(\frac{\mu_{2,1}(1-\mu_{2,0})}{(1-\mu_{2,1})\mu_{2,0}}\right) \\ \vdots \\ \log\left(\frac{\mu_{D,1}(1-\mu_{D,0})}{(1-\mu_{D,1})\mu_{D,0}}\right) \end{bmatrix} \quad (1)$$

$$b = \log\left(\frac{\pi}{1-\pi} \prod_{d=1}^D \frac{1-\mu_{d,1}}{1-\mu_{d,0}}\right) \quad (2)$$

Having found \mathbf{w} and b we can say that we have a linear decision boundary given by hyper-plane equation $\mathbf{w}^T \mathbf{x} + b = 0$