

Week 7:

1) Naive Bayes classifiers:

Let's In Naive Bayes classification, we assume features are independent given class.

$$P(x, y=c) \propto P(x, y=c)$$

$$\begin{aligned} P(y=c/x) &\propto P(x, y=c) \quad (\because \text{Generative model}) \\ &\propto P(x/y=c) P(y=c) \\ &\propto \prod_{i=1}^D P(x_i/y=c) \cdot P(y=c) \end{aligned}$$

MLE for Naive Bayes:

Let's write likelihood for single data case,

$$\begin{aligned} P(x_i, y_i) &= \prod_{j=1}^D P(x_{ij}/y_i) \cdot P(y_i) \\ &= \prod_{j=1}^D \prod_{c=1}^C P(x_{ij}/y_i=c)^{1_{y_i=c}} \cdot \prod_{c=1}^C P(y_i=c)^{1_{y_i=c}} \end{aligned}$$

log likelihood

$$\log P(D) = \sum_{i=1}^N \log \left(\prod_{j=1}^D \prod_{c=1}^C P(x_{ij}/y_i=c)^{1_{y_i=c}} \cdot \prod_{c=1}^C P(y_i=c)^{1_{y_i=c}} \right)$$

$$= \sum_{c=1}^C N_c \log \pi_c + \sum_{j=1}^D \sum_{c=1}^C \sum_{i: y_i=c} P(x_{ij}/y_i=c) \rightarrow (1)$$

[Here we assumed $P(y)$ follows multinomial distribution with $P(y_i=c) = \pi_c$]

and each feature in ' x_j ' can take different kind of distribution. So, for simplicity let's assume, it

follows Bernoulli i.e. $P(x_{ij}/y_i = c) \sim \text{Ber}(\theta_{jc})$

$$\frac{\partial (\log p(\theta))}{\partial \pi_c} = \frac{\partial}{\partial \pi_c} \left(\sum_{i=1}^I \sum_{j=1}^J N_{jc} \log \pi_c \right)$$

Then $\hat{\pi}_c, \hat{\theta}$ can be derived by

$$\underset{\hat{\pi}_c, \hat{\theta}}{\text{argmax}} \log p(\theta) = \underset{\hat{\pi}_c, \hat{\theta}}{\text{argmax}} \text{Expr. (1)}$$

when we are maximizing w.r.t $\hat{\pi}_c, \hat{\theta}$ will be const

$$\text{So } \hat{\pi}_{\text{MLE}} = \underset{\pi}{\text{argmax}} \sum_{c=1}^C N_c \log \pi_c$$

Using MLE derivation for multinomial, this

$$\text{become } \hat{\pi}_{\text{MLE}} = \frac{N_c}{N} = \frac{N_c}{\sum_{i=1}^I N_c}$$

$$\text{and in same way } \hat{\theta}_{jc} = \frac{N_{jc}}{N_c}$$

For Bayesian Model averaging, we take prior for also into account. if we assume prior for π as Dirichlet(α) and for each θ_{jc} as

Beta(β_0, β_1) then prior

$$P(\theta/\theta) P(\theta) = P(\pi) \prod_{j=1}^J \prod_{c=1}^C \pi_c P(\theta_{jc}) \quad (\because \text{factored prior})$$

$$P(\pi) \sim \text{Dir}(\alpha) \quad , \quad P(\theta_{jc}) \sim \text{Beta}(\beta_0, \beta_1)$$

then posterior

$$P(\theta/\theta) = P(\pi/\theta) \prod_{j=1}^J \prod_{c=1}^C \pi_c P(\theta_{jc}/\theta)$$

$P(\pi/\theta)$ is posterior for Dirichlet-Multinomial

which is $p(\pi/\theta) = \text{Dir}(N_1 + \alpha_1, \dots, N_c + \alpha_c)$
 and $p(\theta_{jc}/\theta)$ is posterior of Beta-bernoulli model
 $p(\theta_{jc}/\theta) = \text{Beta}(N_{jc} + \beta_0, (N_{jc} - N_{jc} + \beta_1))$

the Maximum A Posteriori (MAP) estimator for the
 are the mean of the expectation of them (mean)

$$\text{So } \hat{\pi}_{\text{MAP}} = \frac{N_c + \alpha_c}{N + \alpha} \quad (\alpha = \sum_i \alpha_i)$$

$$\text{and } \hat{\theta}_{jc \text{ MAP}} = \frac{N_{jc} + \beta_0}{N_c + \beta_0 + \beta_1}$$

- Here, MAP can be seen as single point estimates for posterior and MLE as single point estimate to explain likelihood
- MAP estimate is also can be seen as weighted average (convex combination) of MLE and prior eq.

13.1: Bayesian inference when σ^2 is unknown:

In linear regression we model $P(y/x) \sim N(w^T x, \sigma^2)$
 In Bayesian modeling, we set priors for w and σ^2 .

$$P(w, \sigma^2) = P(w/\sigma^2) P(\sigma^2)$$

Let's take $P(w/\sigma^2)$ as $N(w_0, \sigma^2 V_0)$ and

$P(\sigma^2)$ as $IG(\sigma^2/a_0, b_0)$.

out prior:

$$P(w, \sigma^2) = N(w/w_0, \sigma^2 V_0) IG(\sigma^2/a_0, b_0)$$

$$= \frac{1}{(2\pi)^{D/2} |\sigma^2 V_0|^{1/2}} \cdot e^{-\frac{1}{2}(w-w_0)^T (\sigma^2 V_0)^{-1} (w-w_0)}$$

$$\times \frac{b_0^{a_0}}{\Gamma(a_0)} \cdot \sigma^{-2(a_0+1)} \cdot e^{-b_0/\sigma^2}$$

$$= \frac{b_0^{a_0}}{(2\pi)^{D/2} |V_0|^{1/2} \Gamma(a_0)} \cdot \sigma^{-2(a_0 + D/2 + 1)} \cdot e^{-\frac{(w-w_0)^T V_0^{-1} (w-w_0) + 2b_0}{2\sigma^2}}$$

likelihood:

$$P(y/x, w, \sigma^2) = \frac{1}{(2\pi)^{D/2}} \cdot N(y/xw, \sigma^2 I_N)$$

$$= \frac{1}{(2\pi)^{D/2} \cdot |\sigma^2 I_N|^{1/2}} \cdot e^{-\frac{1}{2} (y-Xw)^T (\sigma^2 I_N)^{-1} (y-Xw)}$$

$\propto \sigma^{-D/2} \cdot e^{-\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)}$
 posterior for joint distribution of (w, σ^2) is

$$p(w, \sigma^2 / A) = p(y / Xw, \sigma^2) \cdot p(w, \sigma^2)$$

$$\propto \sigma^{-D/2} \cdot e^{-\frac{(y - Xw)^T (y - Xw) + (w - w_0)^T V_0^{-1} (w - w_0)}{2\sigma^2}}$$

which is again in the form of Normal-Inverse Gamma distribution.

$$= \text{NIG}(w, \sigma^2 / w_N, V_N, a_N, b_N)$$

$$w_0 \text{ is updated to } w_N = V_N (V_0^{-1} w_0 + X^T y)$$

$$V_0 \text{ is updated to } V_N = (V_0^{-1} + X^T X)^{-1}$$

$$a_0 \text{ is updated to } a_N = a_0 + N/2$$

$$b_0 \text{ is updated to } b_N = b_0 + \frac{1}{2} (w_0^T V_0^{-1} w_0 + y^T y - w_N^T V_N^{-1} w_N)$$

4.19) Generative classifier

$$P(y=c|x) = \frac{P_0(y=c) P_0(x/y=c)}{\sum_{c'} P_0(y=c') P_0(x/y=c')}$$

In Gaussian discriminant, we model

$$P(y) = \text{MNL}(\pi) \quad P(x/y=c) = N(\mu_c, \Sigma_c)$$

$$\theta = [\pi_{1:c}, \mu_{1:c}, \Sigma_{1:c}]$$

In given problem, we have only two classes $c=2$

and $\Sigma_1 = K \Sigma_0$, $K > 0$, then we can write

$$P(y=1/x) = \frac{\pi_1 \cdot |2\pi \Sigma_1|^{-1/2} e^{-\frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1)}}{\sum_{i=0}^1 \pi_i \cdot |2\pi \Sigma_i|^{-1/2} e^{-\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1}(x-\mu_i)}}$$

We'll simplify numerator & denominator separately

by substituting $K \Sigma_0 = \Sigma_1$, $K > 0$

$$\text{Numerator of } P(y=1/x) = \pi_1 \cdot |2\pi (K \Sigma_0)|^{-1/2} e^{-\frac{1}{2}(x-\mu_1)^T (K \Sigma_0)^{-1}(x-\mu_1)}$$

$$= \frac{1}{\sqrt{K}} \cdot \pi_1 \cdot |2\pi \Sigma_0|^{-1/2} \cdot e^{-\frac{1}{2K}(x-\mu_1)^T \Sigma_0^{-1}(x-\mu_1)}$$

Denominator of $P(y=1/x) =$

$$= \pi_1 \cdot |2\pi (K \Sigma_0)|^{-1/2} e^{-\frac{1}{2}(x-\mu_1)^T (K \Sigma_0)^{-1}(x-\mu_1)} + \pi_0 \cdot |2\pi \Sigma_0|^{-1/2} e^{-\frac{1}{2}(x-\mu_0)^T \Sigma_0^{-1}(x-\mu_0)}$$

$$= \frac{1}{\sqrt{k}} (\pi_1 \cdot (2\pi \Sigma_0)^{-1})$$

$$= \frac{1}{\sqrt{k}} (\pi_1 \cdot (2\pi \Sigma_0)^{-1}) \cdot e^{-\frac{1}{2} x^T \Sigma_0^{-1} x} \left[\frac{\pi_1}{\sqrt{k}} \cdot e^{\frac{1}{k} (\mu_1^T \Sigma_0^{-1} x - \frac{1}{2} \mu_1^T \Sigma_0 \mu_1)} + \frac{\pi_0}{\sqrt{k}} \cdot e^{\frac{1}{k} (\mu_0^T \Sigma_0^{-1} x - \frac{1}{2} \mu_0^T \Sigma_0 \mu_0)} \right]$$

After we cancel first term i.e. $\frac{1}{\sqrt{k}} (\pi_1 \cdot (2\pi \Sigma_0)^{-1}) \cdot e^{-\frac{1}{2} x^T \Sigma_0^{-1} x}$

both in numerator & denominator we get.

$$= \frac{\left(\frac{\pi_1}{\sqrt{k}} \cdot e^{\frac{1}{k} (\mu_1^T \Sigma_0^{-1} x - \frac{1}{2} \mu_1^T \Sigma_0 \mu_1)} \right)}{\left(\frac{\pi_1}{\sqrt{k}} \cdot e^{\frac{1}{k} (\mu_1^T \Sigma_0^{-1} x - \frac{1}{2} \mu_1^T \Sigma_0 \mu_1)} + \frac{\pi_0}{\sqrt{k}} \cdot e^{\frac{1}{k} (\mu_0^T \Sigma_0^{-1} x - \frac{1}{2} \mu_0^T \Sigma_0 \mu_0)} \right)}$$

$$\frac{\pi_1}{\sqrt{k}} \cdot e^{\frac{1}{k} (\mu_1^T \Sigma_0^{-1} x - \frac{1}{2} \mu_1^T \Sigma_0 \mu_1)} + \frac{\pi_0}{\sqrt{k}} \cdot e^{\frac{1}{k} (\mu_0^T \Sigma_0^{-1} x - \frac{1}{2} \mu_0^T \Sigma_0 \mu_0)}$$

Let's simplify this by substituting

$$\beta_1 = \Sigma_0^{-1} \mu_1$$

$$\beta_0 = \Sigma_0^{-1} \mu_0$$

$$r_1 = -\frac{1}{2} \mu_1^T \Sigma_0 \mu_1 + \log \pi_1 \quad r_0 = -\frac{1}{2} \mu_0^T \Sigma_0 \mu_0 + \log \pi_0$$

$$= \frac{e^{\frac{1}{k} (\beta_1^T x + r_1 - \frac{1}{2} \log k)}}{e^{\frac{1}{k} (\beta_1^T x + r_1 - \frac{1}{2} \log k)} + e^{\frac{1}{k} (\beta_0^T x + r_0 - \frac{1}{2} \log k)}}$$

$$= \frac{1}{1 + e^{\beta_0^T x + r_0 - \frac{1}{k} (\beta_1^T x + r_1 + \frac{1}{2} \log k)}}$$

$$= \frac{1}{1 + e^{(\beta_0 - \gamma_k \beta_1)^T x + (r_0 - \gamma_k r_1) + \frac{\log k}{2k}}}$$

$$= \text{Sigm} \left((\beta_{1/k} - \beta_0)^T x + (\gamma_{1/k} r_1 - r_0) + \frac{\log k}{2k} \right)$$

So, Even $\Sigma_1 = k \Sigma_0$ the discriminative ~~surface~~ function is still sigmoid function. and boundary is still linear.