

10-Generative Models

Discriminative models vs. Generative models

- In the training of discriminative models, we have a generative story for t given x .
For example,

- Linear regression: $t|x \sim N(w^T x + b, \sigma_e^2)$
- Softmax classification: $t|x \sim \text{Multinomial}(\text{softmax}(Wx + b))$

And we model the conditional probability of t given x .

$$p(t|x; \Theta)$$

In discriminative models, we do NOT have a generative story of x . In other words, we do not have a model for $p(x)$.

- In **generative models**, we will have a generative story concerning how both x and t are generated, i.e., we will have a probabilistic model of $p(x, t)$.

- First generate x , and then generate t given x .

$$p(x, t) = p(x)p(t|x)$$

x generated is more complicated than y , but we model x iid. We then model $p(t|x)$, which is not different from the discriminative model.

=> Not quite meaningful

- First generate t and then generate x given t .

$$p(x, t) = p(t)p(x|t)$$

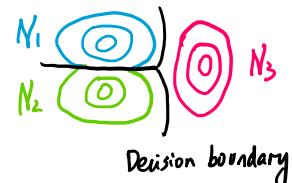
For classification, y is a categorical variable. $p(t)$ models class frequency, and $p(x|t)$ models the distribution of x in a certain category.

Since we will model how x is obtained given t , we need to distinguish the type of x , i.e., continuous variables or discrete variables.

Gaussian discriminant analysis

- Generative story for discrete y and continuous features x .

$$\begin{aligned} t &\sim \text{Multinomial}(\pi_1, \dots, \pi_K) \\ x|t = k &\sim N(\mu_k, \Sigma_k) \end{aligned}$$



- Parameter estimation: Maximum likelihood estimation

$$\begin{aligned}
 \hat{\pi}, \hat{\mu}_1, \dots, \hat{\mu}_K, \hat{\Sigma}_1, \dots, \hat{\Sigma}_K &= \operatorname{argmax} \log p(\mathbf{x}, \mathbf{t}) \quad [\text{Generative model maximizes joint likelihood}] \\
 &= \operatorname{argmax} \sum_{m=1}^M \log p(\mathbf{x}^{(m)}, t^{(m)}) \quad [\text{data iid}] \\
 &= \operatorname{argmax} \sum_{m=1}^M \log [p(t^{(m)}) p(\mathbf{x}^{(m)} | t^{(m)})] \quad [\text{Generative story}] \\
 &= \operatorname{argmax} \sum_{m=1}^M [\log p(t^{(m)}) + \log p(\mathbf{x}^{(m)} | t^{(m)})] \\
 &= \operatorname{argmax} \left[\sum_{m=1}^M \log p(t^{(m)}; \pi) + \sum_{m=1}^M \log p(\mathbf{x}^{(m)} | t^{(m)}; \mu_{t^{(m)}}, \Sigma_{t^{(m)}}) \right] \\
 &= \operatorname{argmax} \left[\sum_{m=1}^M \log p(t^{(m)}; \pi) + \sum_{k=1}^K \sum_{\substack{m=1..M \\ t^{(m)}=k}} p(\mathbf{x}^{(m)} | t^{(m)}=k; \mu_k, \Sigma_k) \right] \\
 &\qquad\qquad\qquad [\text{rearrange of terms}]
 \end{aligned}$$

This shows that all parameters decompose into the following groups:

$$\pi, (\mu_1, \Sigma_1), \dots, (\mu_K, \Sigma_K)$$

We may optimize each group separately.

- MLE for multinomial distribution

- Quick result (just counting):

$$\hat{\pi}_k = \frac{\sum_{m=1}^M \mathbb{1}\{t^{(m)}=k\}}{M} \quad \text{for } k=1, \dots, K$$

- 1st Proof: Set $\bar{\pi}_K = 1 - \pi_1 - \dots - \pi_{K-1}$ and seek for a closed form solution
- 2nd Proof (Lagrange multiplier method):

$$\begin{aligned}
 \log p(t^{(1)}, \dots, t^{(M)}) &= \sum_{m=1}^M \log p(t^{(m)}) \\
 &= \sum_{m=1}^M \log \bar{\pi}_{t^{(m)}} \\
 &= \sum_{m=1}^M \sum_{k=1}^K \mathbb{1}\{t^{(m)}=k\} \log \pi_k
 \end{aligned}$$

The optimization problem becomes

$$\text{minimize} \quad - \sum_{m=1}^M \sum_{k=1}^K \mathbb{1}\{t^{(m)}=k\} \log \pi_k$$

$$\text{subject to} \quad \pi_1 + \pi_2 + \dots + \pi_K = 1$$

$$\text{Lagrangian: } \mathcal{L}(\pi, \lambda) = - \sum_{m=1}^M \sum_{k=1}^K [\mathbb{1}\{t^{(m)}=k\} \cdot \log \pi_k] + \lambda(\pi_1 + \dots + \pi_K - 1)$$

Optimality conditions:

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}(\pi, \lambda)}{\partial \pi_k} = 0 \\ \frac{\partial \mathcal{L}(\pi, \lambda)}{\partial \lambda} = 0 \end{array} \right. \quad \begin{array}{l} (1) \\ (2) \end{array}$$

(1) implies

$$\frac{\partial \mathcal{L}(\pi, \lambda)}{\partial \pi_k} = - \sum_{m=1}^M \mathbb{1}\{t^{(m)}=k\} \cdot \frac{1}{\pi_k} + \lambda = 0 \quad (3)$$

(2) implies

$$\pi_1 + \dots + \pi_K = 1 \quad (4)$$

(3) implies

$$\lambda \pi_k = \sum_{m=1}^M \mathbb{1}\{t^{(m)}=k\} \quad (5)$$

Put (5) in (4):

$$\sum_{k=1}^K \sum_{m=1}^M \mathbb{1}\{t^{(m)}=k\} = \lambda$$

$$\text{which implies } \lambda = M \quad (6)$$

Put (6) back in (5):

$$\hat{\pi}_k = \frac{\sum_{m=1}^M \mathbb{1}\{t^{(m)}=k\}}{M}$$

- MLE for Gaussian distributions

$$\hat{\mu}_k = \frac{1}{M_k} \sum_{m=1}^M \mathbb{1}\{t^{(m)}=k\} \mathbf{x}_k$$

$$\hat{\Sigma}_k = \frac{1}{M_k} \sum_{m=1}^M \mathbb{1}\{t^{(m)}=k\} \cdot (\mathbf{x}^{(m)} - \hat{\mu}_k)(\mathbf{x}^{(m)} - \hat{\mu}_k)^T$$

$$\text{where } M_k = \sum_{m=1}^M \mathbb{1}\{t^{(m)}=k\}$$

If we have an assumption that all covariance $\Sigma_1, \dots, \Sigma_K$ are shared as Σ , then

$$\hat{\Sigma} = \frac{M_1}{M} \hat{\Sigma}_1 + \dots + \frac{M_K}{M} \hat{\Sigma}_K$$

- Decision boundary for Gaussian discriminant analysis

Max *a posteriori* inference: $\hat{t} = \operatorname{argmax} \hat{p}(t | \mathbf{x}) = \operatorname{argmax} \hat{p}(t) p(\mathbf{x} | t)$

Considering any two categories i, j .

$\hat{t} = i$ being preferred than $\hat{t} = j$

$$\Leftrightarrow \hat{p}(t=i|x) \geq \hat{p}(t=j|x)$$

$$\Leftrightarrow \pi_i N(x; \mu_i, \Sigma_i) \geq \pi_j N(x; \mu_j, \Sigma_j)$$

$$\Leftrightarrow \pi_i \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right\}$$

$$\geq \pi_j \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j) \right\}$$

\Leftrightarrow

$$-\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) + \log \frac{\pi_i}{|\Sigma_i|^{1/2}} \geq -\frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j) + \log \frac{\pi_j}{|\Sigma_j|^{1/2}}$$

If $\Sigma_i \neq \Sigma_j$, then the decision boundary is quadratic.

If $\Sigma_i = \Sigma_j$, then the decision boundary is linear

[See the above RGB cartoon]

Naïve Bayes for discrete variables

- For simplicity, we only consider binary features
 $x_i \in \{0,1\}$, i.e., $x \in \{0,1\}^d$
- Generative story for discrete y and (binary) discrete features x
 $t \sim \text{Multinomial}(\pi_1, \dots, \pi_K)$
 $x_i | t = k \sim \text{Bernoulli}(p_{k,i})$
 - Such model can be used to represent a document in text classification.
For example, $t = 1$ indicates Spam, $t = 0$ indicates NotSpam. $x_i \in \{0,1\}$ indicates if the i th word in the vocabulary occurs in the document.
- MLE for Naïve Bayes [HW]
Please show that the parameters of naïve Bayes decompose, i.e., the probability factorizes (for the same reason as Gaussian discriminative analysis). Write out the MLE for naïve Bayes (which is simply counting).

*Naïve Bayes is named for the assumption that the MLE for

NO proof is needed for the second part, because the rule for multinomial distribution has been clear in the Gaussian discriminative analysis.

- Inference with the naïve Bayes model

Category k is preferred than Category j

$$\begin{aligned} \Leftrightarrow \pi_k \prod_{i=1}^d p_{k,i}^{x_i} (1-p_{k,i})^{1-x_i} &\geq \pi_j \prod_{i=1}^d p_{j,i}^{x_i} (1-p_{j,i})^{1-x_i} \\ \Leftrightarrow \log \pi_k + \sum_{i=1}^d [x_i \log p_{k,i} + (1-x_i) \log (1-p_{k,i})] &\geq \log \pi_j + \sum_{i=1}^d [x_i \cdot \log p_{j,i} + (1-x_i) \log (1-p_{j,i})] \end{aligned}$$

The decision boundary is again linear.

Comparing discriminative models and generative models

- For naïve Bayes and Gaussian discriminative analysis with shared covariance, the decision boundary is linear, which yield the same hypothesis class as logistic regression.

- Even for Gaussian discriminative analysis with different covariance matrices, the decision boundary is quadratic. We may introduce quadratic features for logistic regression, too.

Thus, generative models do not enhance the hypothesis class, compared with discriminative models.

- Generative and discriminative models have the same inference criterion (when we concern the accuracy as our measure of success).

$$\hat{t}_* = \operatorname{argmax} \hat{p}(t|\mathbf{x})$$

However, they differ in the training objective. Discriminative models treat input \mathbf{X} as constants (unknown variables), whereas generative models model \mathbf{X} as random variables.

$$\hat{\theta} = \operatorname{argmax} p(\mathbf{y}|\mathbf{X}; \theta) \text{ vs. } \hat{\theta} = \operatorname{argmax} p(\mathbf{X}, \mathbf{y}; \theta)$$

Note: MAP parameter estimation can be applied to both.

- Generative models usually have far more parameters than discriminative models, thus requiring more data for training.
- Generative models work for more scenarios than classification.

- Unsupervised learning. MLE is still well defined in generative models if \mathbf{X} is not given. [Additional lecture]
- Inferring some of x_i during prediction, given other features and labels.