

Generative Models for Supervised Learning,

CS772A: Probabilistic Machine Learning

Piyush Rai

Generative Supervised Learning

- The conditional distribution $p(y|x)$ can also be defined as

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

Requires modeling the joint distribution of the inputs and outputs

In the discriminative approach for learning $p(y|x)$, we didn't model the inputs x but treated them as "given"

- Generative sup. learning is usually more work because $p(x, y)$ has to be estimated
- However, there are some benefits as well. For example, for classification

$p(y)$ is called the "class-prior" or "class-marginal" distribution

Can incorporate knowledge of frequency ("size") of each class in training data

Can incorporate knowledge of the distribution ("shape") of each class in training data

$$p(y|x) = \frac{p(x, y)}{p(x)} = \frac{p(y)p(x|y)}{p(x)}$$

Can assume simple/sophisticated types of distributions for the "class-conditional" distribution $p(x|y)$ and learned them using the training data of each class

- Once class-marginal and class-conditionals are estimated, we can make predictions for the label of test input as $p(y_* = k|x_*) \propto p(y_* = k)p(x_*|y_* = k)$



Generative Supervised Learning



- The generative classification model

Probability of belonging to class k , conditioned on the input \mathbf{x}

$$p(y = k|\mathbf{x}) = \frac{p(y = k)p(\mathbf{x}|y = k)}{\sum_k p(y = k)p(\mathbf{x}|y = k)}$$

Marginal probability of belonging to class k

Probability (density) of input \mathbf{x} under class k

Note: Estimating $p(\mathbf{x}|\mathbf{y})$ can be difficult especially if \mathbf{x} is high-dimensional and we don't have enough data from each class

A way to handle this is to assume **simpler forms** for $p(\mathbf{x}|\mathbf{y})$ (e.g., Gaussian with diagonal/spherical covar – **naïve Bayes**) but it might sacrifice accuracy too

- We need to learn $p(\mathbf{y})$ and $p(\mathbf{x}|\mathbf{y})$ here given training data $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$

- Class prior/marginal distribution $p(\mathbf{y})$ will always be a discrete distribution, e.g.,

- For $y \in \{0,1\}$, $p(y) = p(y|\pi) = \text{Bernoulli}(y|\pi)$ with $\pi \in (0,1)$

- For $y \in \{1,2,\dots,K\}$, $p(y) = p(y|\pi) = \text{multinoulli}(y|\pi)$ where $\pi = [\pi_1, \dots, \pi_K]$

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

- Class conditional distribution $p(\mathbf{x}|\mathbf{y})$ will depend on the nature of inputs, e.g.,

- For $\mathbf{x} \in \mathbb{R}^D$, $p(\mathbf{x}|\mathbf{y} = k)$ can be a multivariate Gaussian (one per class)

For π , can use Beta or Dirichlet (we have already seen these examples)

Note: When estimating θ_k , we only need inputs from class k
 $\mathbf{X}_k = \{\mathbf{x}_n: y_n = k\}$

$$p(\mathbf{x}|\mathbf{y} = k) = p(\mathbf{x}|\theta_k) = \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)$$

Will need appropriate prior distributions for π and $\{\theta_k\}_{k=1}^K$

- Can estimate π and $\{\theta_k\}_{k=1}^K$ using (\mathbf{X}, \mathbf{y}) via point est. or fully Bayesian infer.

Generative Classification: Making Predictions

- Once π and $\{\theta_k\}_{k=1}^K$ are learned, we are ready to make prediction for any test input \mathbf{x}_*
- Two ways to make the prediction
- Approach 1: If we have point estimates for π and $\{\theta_k\}_{k=1}^K$, say $\hat{\pi}$ and $\{\hat{\theta}_k\}_{k=1}^K$. Then

$$p(y_* = k | \mathbf{x}_*) = \frac{p(y_* = k | \hat{\pi}) p(\mathbf{x}_* | \hat{\theta}_k)}{\sum_k p(y_* = k | \hat{\pi}) p(\mathbf{x}_* | \hat{\theta}_k)} \propto \hat{\pi}_k p(\mathbf{x}_* | \hat{\theta}_k)$$

Compute for every value of k and normalize

- Approach 2: If we have the full posterior for π and $\{\theta_k\}_{k=1}^K$. Then
 - Instead of using $p(y_* = k | \hat{\pi})$, we will use $p(y_* = k | \mathbf{y}) = \int p(y_* = k | \pi) p(\pi | \mathbf{y}) d\pi$ PPD of y_*
 - Instead of using $p(\mathbf{x}_* | \hat{\theta}_k)$, we will use $p(\mathbf{x}_* | \mathbf{X}_k) = \int p(\mathbf{x}_* | \theta_k) p(\theta_k | \mathbf{X}_k) d\theta_k$ PPD of \mathbf{x}_*
 - Using these quantities, the prediction will be made as

$$p(y_* = k | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \frac{p(y_* = k | \mathbf{y}) p(\mathbf{x}_* | \mathbf{X}_k)}{\sum_k p(y_* = k | \mathbf{y}) p(\mathbf{x}_* | \mathbf{X}_k)} \propto p(y_* = k | \mathbf{y}) p(\mathbf{x}_* | \mathbf{X}_k)$$

Compute for every value of k and normalize

Note that we aren't using a single "best" value of the params π and θ_k unlike Approach 1



Generative Classification: A Generative Story

- Assuming binary labels, can define a “generative story” for each example (\mathbf{x}_i, y_i)
 - First draw (“generate”) a binary label $y_i \in \{0,1\}$

$$y_i \sim \text{Bernoulli}(\pi)$$

For multi-class problems, we will have a multinoulli instead

- Now draw (“generate”) the input \mathbf{x}_i from the distribution of class $y_i \in \{0,1\}$

$$\mathbf{x}_i | y_i \sim p(\mathbf{x} | \theta_{y_i})$$

Most generative models (supervised as well as unsupervised/semi-supervised) can be expressed via such a story



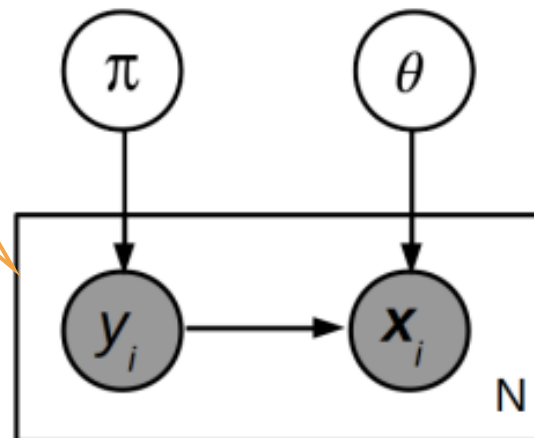
- Writing $\theta = (\theta_0, \theta_1)$, the above generative model shown in “plate notation”

Order of generation in this story depends on what part of the data/parameters depend on what data/params

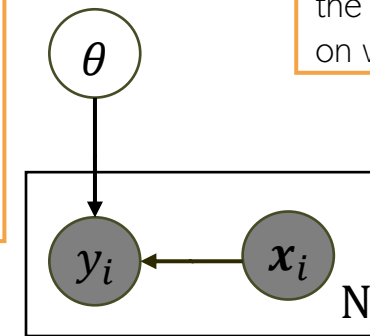
Often we also show the generation of parameters/unknowns as well (via their respective distributions)

Note that in this generative process, we assume y is generated first since the generation of \mathbf{x} depends on what y is

Note that in this generative process, we assume y is generated first since the generation of \mathbf{x} depends on what y is



Compare this with the plate notation diagram of a discriminative model such as prob linear regression or logistic regression



A discriminative model (no model for \mathbf{x}_i 's)

Generative Sup. Learning: Some Comments

- A very flexible approach for classification

Incorporate info about how frequent each class is in the training data ("class prior")

Incorporate info about the shape of each class

Consequently, can naturally learn nonlinear boundaries, too (without using kernel methods or deep learning)

$$p(y_* = k | \mathbf{x}_*) = \frac{p(y_* = k)p(\mathbf{x}_* | y_* = k)}{\sum_k p(y_* = k)p(\mathbf{x}_* | y_* = k)}$$

- Can handle **missing labels** and **missing features**

Will discuss this later

- These can be treated as latent variables as estimated using methods such as EM

- Ability to handle missing labels makes it suitable for **semi-supervised learning**

- The choice of the class-conditional and proper estimation is important

- Can leverage advances in deep generative models to learn very flexible forms for $p(\mathbf{x}|y)$

- Can also use it for **regression** (define $p(\mathbf{x}, y)$ via some distr. and obtain $p(y|\mathbf{x})$ from it)

- Can also combine generative and discriminative approaches for supervised learning



Hybrids of Discriminative and Generative Models

7

- Both discriminative and generative models have their strengths/shortcomings
- Some aspects about discriminative models for sup. learning
 - Discriminative models have usually fewer parameters (e.g., just a weight vector)
 - Given “plenty” of training data, disc. models can usually outperform generative models
- Some aspects about generative models for sup. learning
 - Can be more flexible (we have seen the reasons already)
 - Usually have more parameters to be learned
 - Modeling the inputs (learning $p(\mathbf{x}|\mathbf{y})$) can be difficult for high-dim inputs
- Some prior work on combining discriminative and generative models. Examples:

Recall prob linear regression and logistic reg

$$\alpha \log p(y|x; \theta) + \beta \log p(x; \theta)$$

Approach 1 (McCullum et al, 2006) – modeling the joint $p(\mathbf{x}, \mathbf{y}|\theta)$ using a multi-conditional likelihood

$$p(x, y, \theta_d, \theta_g) = p_{\theta_d}(y|x)p_{\theta_g}(x)p(\theta_d, \theta_g)$$

Approach 2 (Lasserre et al, 2006) – Coupled parameters between discriminative and generative models

$$p(x, y, z) = p(y|x, z) \cdot p(x, z)$$

Approach 3 (Kuleshov and Ermon, 2017) – Coupling discriminative and generative models via a latent variable \mathbf{z} (see “Deep Hybrid Models: Bridging Discriminative and Generative Approaches”, UAI 2017)