

Exponential Family Distributions (contd.), Generative Models for Supervised Learning

CS772A: Probabilistic Machine Learning

Piyush Rai

Plan today

- Exponential Family Distributions
 - Conjugate priors, posterior, and PPD
- Generative approach to supervised learning



Bayesian Inference for Expon. Family Distributions³

- Already saw that the total **likelihood** given N i.i.d. observations $\mathcal{D} = \{x_1, \dots, x_N\}$

$$p(\mathcal{D}|\theta) \propto \exp \left[\theta^\top \phi(\mathcal{D}) - NA(\theta) \right] \quad \text{where} \quad \phi(\mathcal{D}) = \sum_{i=1}^N \phi(x_i)$$

- Let's choose the following **prior** (note: looks similar in terms of θ within exp)

$$p(\theta|\nu_0, \tau_0) = h(\theta) \exp \left[\theta^\top \tau_0 - \nu_0 A(\theta) - A_c(\nu_0, \tau_0) \right]$$

- Ignoring the prior's log-partition function $A_c(\nu_0, \tau_0) = \log \int_{\theta} h(\theta) \exp \left[\theta^\top \tau_0 - \nu_0 A(\theta) \right] d\theta$

$$p(\theta|\nu_0, \tau_0) \propto h(\theta) \exp \left[\theta^\top \tau_0 - \nu_0 A(\theta) \right]$$

- Comparing the prior's form with the likelihood, note that
 - ν_0 is like the number of “pseudo-observations” coming from the prior
 - τ_0 is the total sufficient statistics of the pseudo-observations (τ_0 / ν_0 per pseudo-obs)



The Posterior

- The likelihood and prior were

$$p(\mathcal{D}|\theta) \propto \exp \left[\theta^\top \phi(\mathcal{D}) - NA(\theta) \right] \quad \text{where} \quad \phi(\mathcal{D}) = \sum_{i=1}^N \phi(\mathbf{x}_i)$$

Assume its log partition function denoted as $A_c(\nu_0, \tau_0)$

$$p(\theta|\nu_0, \tau_0) \propto h(\theta) \exp \left[\theta^\top \tau_0 - \nu_0 A(\theta) \right]$$

Posterior is also from the same family as the prior

Happens when the prior is conjugate to the likelihood

- The posterior $p(\theta|\mathcal{D}) \propto p(\theta)p(\mathcal{D}|\theta)$ therefore will be

Its log partition function will be $A_c(\nu_0 + N, \tau_0 + \phi(\mathcal{D}))$

$$p(\theta|\mathcal{D}) \propto h(\theta) \exp \left[\theta^\top (\tau_0 + \phi(\mathcal{D})) - (\nu_0 + N)A(\theta) \right]$$

- Every exp family likelihood has a conjugate prior having the form above
- Posterior's hyperparams τ'_0, ν'_0 obtained by adding "stuff" to prior's hyperparams

Number of pseudo-observations plus number of actual observations

$$\nu'_0 \leftarrow \nu_0 + N$$

Suff-stats of pseudo-observations plus suff-stats of actual observations

$$\tau'_0 \leftarrow \tau_0 + \phi(\mathcal{D})$$

Another equivalent form

$$p(\theta|\mathcal{D}) \propto h(\theta) \exp \left[\theta^\top (\nu_0 + N) \frac{\nu_0 \bar{\tau}_0 + \phi(\mathcal{D})}{\nu_0 + N} - (\nu_0 + N)A(\theta) \right]$$

$$\bar{\tau}_0 = \tau_0 / \nu_0$$

$$\begin{aligned} \nu'_0 &\leftarrow \nu_0 + N \\ \bar{\tau}'_0 &\leftarrow \frac{\nu_0 \bar{\tau}_0 + N \bar{\phi}}{\nu_0 + N} \end{aligned}$$

$$\bar{\phi} = \frac{\phi(\mathcal{D})}{N}$$

Convex comb of avg suff-stats of pseudo obs and actual obs



Posterior Predictive Distribution

- Assume some training data $\mathcal{D} = \{x_1, \dots, x_N\}$ from some exp-fam distribution
- Assume some test data $\mathcal{D}' = \{\tilde{x}_1, \dots, \tilde{x}_{N'}\}$ from the same distribution
- The posterior pred. distr. of \mathcal{D}'

$$\begin{aligned}
 p(\mathcal{D}'|\mathcal{D}) &= \int \underbrace{p(\mathcal{D}'|\theta)}_{\text{Exp. Fam. likelihood w.r.t. test data}} \underbrace{p(\theta|\mathcal{D})}_{\text{Posterior (same form as the prior due to conjugacy)}} d\theta \\
 &= \int \underbrace{\left[\prod_{i=1}^{N'} h(\tilde{x}_i) \right]}_{\text{constant w.r.t. } \theta} \exp \left[\theta^\top \phi(\mathcal{D}') - N' A(\theta) \right] h(\theta) \exp \left[\theta^\top (\tau_0 + \phi(\mathcal{D})) - (\nu_0 + N) A(\theta) - \underbrace{A_c(\nu_0 + N, \tau_0 + \phi(\mathcal{D}))}_{\text{constant w.r.t. } \theta} \right] d\theta
 \end{aligned}$$

- This gets further simplified into

$$\begin{aligned}
 p(\mathcal{D}'|\mathcal{D}) &= \left[\prod_{i=1}^{N'} h(\tilde{x}_i) \right] \frac{\int h(\theta) \exp \left[\theta^\top (\tau_0 + \phi(\mathcal{D}) + \phi(\mathcal{D}')) - (\nu_0 + N + N') A(\theta) \right] d\theta}{\exp [A_c(\nu_0 + N, \tau_0 + \phi(\mathcal{D}))]} \\
 &= \left[\prod_{i=1}^{N'} h(\tilde{x}_i) \right] \frac{Z_c(\nu_0 + N + N', \tau_0 + \phi(\mathcal{D}) + \phi(\mathcal{D}'))}{\exp [A_c(\nu_0 + N, \tau_0 + \phi(\mathcal{D}))]}
 \end{aligned}$$



Posterior Predictive Distribution

- Since $A_c = \log Z_c$ or $Z_c = \exp(A_c)$, we can write the PPD as

$$\begin{aligned} p(\mathcal{D}'|\mathcal{D}) &= \left[\prod_{i=1}^{N'} h(\tilde{\mathbf{x}}_i) \right] \frac{Z_c(\boldsymbol{\nu}_0 + N + N', \boldsymbol{\tau}_0 + \phi(\mathcal{D}) + \phi(\mathcal{D}'))}{Z_c(\boldsymbol{\nu}_0 + N, \boldsymbol{\tau}_0 + \phi(\mathcal{D}))} \\ &= \left[\prod_{i=1}^{N'} h(\tilde{\mathbf{x}}_i) \right] \exp [A_c(\boldsymbol{\nu}_0 + N + N', \boldsymbol{\tau}_0 + \phi(\mathcal{D}) + \phi(\mathcal{D}')) - A_c(\boldsymbol{\nu}_0 + N, \boldsymbol{\tau}_0 + \phi(\mathcal{D}))] \end{aligned}$$

Thus PPD as well as marginal likelihood has closed form expression when working with exp-family distributions



- Therefore the **posterior predictive** is proportional to
 - Ratio of two partition functions of two “posterior distributions” (one with $N + N'$ examples and the other with N examples)
 - Exponential of the difference of the corresponding log-partition functions
- Note that the form of Z_c (and A_c) will simply depend on the chosen conjugate prior
- Very useful result. Also holds for $N = 0$
 - In this case $p(\mathcal{D}') = \int p(\mathcal{D}'|\theta)p(\theta)d\theta$ is simply the **marginal likelihood** of test data \mathcal{D}'



Summary

- Exp. family distributions are very useful for modeling diverse types of data/parameters
- Conjugate priors to exp. family distributions make parameter updates very simple
- Other quantities such as posterior predictive can be computed in closed form
- Useful in designing generative classification models. Choosing class-conditional from exponential family with conjugate priors helps in parameter estimation
- Useful in designing generative models for unsupervised learning
- Used in designing **Generalized Linear Models**: Model $p(\mathbf{y}|\mathbf{x})$ using exp. fam distribution
 - Linear regression (with Gaussian likelihood) and logistic regression are GLMs
- Will see several use cases when we discuss approx inference algorithms (e.g., Gibbs sampling, and especially variational inference)



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



Generative Supervised Learning



- The generative classification model

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

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

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

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

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

- 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 $\boldsymbol{\pi}$, 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 $\boldsymbol{\pi}$ and $\{\theta_k\}_{k=1}^K$

- Can estimate $\boldsymbol{\pi}$ 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 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 ¹²

- 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)