#### Probabilistic ML: Some Basic Ideas

CS772A: Probabilistic Machine Learning Piyush Rai

### Probabilistic ML Modeling: The Basic Ingredients

• Likelihood model  $p(\mathcal{D}|\theta)$  for data  $\mathcal{D}$ ; prior distribution  $p(\theta|\alpha)$  over parameters  $\theta$ 



- Likelihood defined in terms of distribution(s) we assume data is generated from
  - It's like a measure of "fit" between observed data and each possible value of parameters
  - Its negative is like the "loss function" (high likelihood value = low loss; and vice-versa)
- Prior specifies our prior knowledge about  $\theta$  before we have seen the data
  - It also acts as a regularizer for  $\theta$  (will see the reason formally later)
- Note: The prior itself depends on other parameters lpha (also unknown)
  - These are sometimes called "hyperparameters" (can set by hand or estimate from data): PML

#### The Prior: Where does it come from?

- The prior  $p(\theta | \alpha)$  plays an important role in probabilistic/Bayesian modeling
  - Reflects our prior beliefs about possible parameter values <u>before</u> seeing the data



- Can be "subjective" or "objective" (also a topic of debate, which we won't get into)
- Subjective: Prior (our beliefs) derived from past experiments
- Objective: Prior represents "neutral knowledge" (e.g., uniform, vague prior)
- Can also be seen as a regularizer (connection with non-probabilistic view)



#### Parameter Estimation

- lacksquare The parameters m heta are unknown and need to be estimated from training data  $\mathcal D$
- When estimating  $\theta$ , we may take one of the following approaches Approach 1 Approach 2
  - heta has an unknown with fixed value
  - Estimate the single best estimate of  $\theta$  by optimizing a loss function  $\hat{\theta} = \operatorname{argmin}_{\theta} \mathcal{L}(\mathcal{D}; \theta)$

Treat  $\theta$  as a random variable

• Estimate  $\theta$  by computing its distribution conditioned on  $\mathcal{D}$ Posterior distribution  $p(\theta | \mathcal{D})$ 

• Approach 2 also gives uncertainty about our estimate of  $\theta$ ; Approach 1 doesn't

• But possible to estimate uncertainty in  $\theta$  even with Approach 1 (e.g., using ensembles).

- Approach 1 is also a simplified/special case/approximation of Approach 2
- Can also take a hybrid (Approach 2 for some parameters; Approach 1 for others).

#### The Posterior Distribution

The posterior distribution is computed using Bayes rule (Bayesian inference)



#### "Online" Nature of Bayesian Inference Updates

Bayesian inference can naturally be done in an online fashion



CS772A: PML

#### Point Estimation

Recall that the posterior is

Intractable to compute except for some very simple models or if the likelihood and prior are **conjugate** (discussed later) to each other Intractable mainly because the marginal likelihood (the denominator on the RHS is intractable in general)

 $p(\mathcal{D}|\theta)p(\theta|\alpha)$ 

However, point estimation throws away all the uncertainty information about  $\theta$ 

Meaning the observed data has the

largest probability for this value of heta

CS772A: PML

If posterior is intractable, can use MLE/MAP to get point estimates

• Maximum likelihood (ML) estimation: Find  $\theta$  for which likelihood is highest

 $p(\theta | \mathcal{D}, \alpha) =$ 

Negative Log likelihood (equivalent to a loss function)

$$\hat{\theta}_{ML} = \arg\max_{\theta} \log \frac{p(\mathcal{D}|\theta)}{\theta} = \arg\min_{\theta} - \log \frac{p(\mathcal{D}|\theta)}{\theta} = \arg\min_{\theta} NLL(\theta)$$

• Maximum a posteriori (MAP) estimation: Find  $\theta$  with largest posterior prob.

$$\hat{\theta}_{MAP} = \arg\max_{\theta} \log p(\theta | \mathcal{D}, \alpha) = \arg\max_{\theta} \left[ \log p(\mathcal{D} | \theta) + \log p(\theta | \alpha) \right]$$
  
Like MLE with info from prior added =  $\arg\min_{\theta} \left[ NLL(\theta) - \log p(\theta | \alpha) \right]$ 

Akin to a regularizer added to the loss is part of prior

#### The Predictive Distribution

- Predictive distribution is the distribution of test data  $\mathcal{D}_*$  given training data  $\mathcal{D}$
- In the general form, we can write it as



• If we only have point estimate of  $\theta$  (say  $\hat{\theta}$  obtained from MLE/MAP) then

This approximation of PPD is called "plug-in" predictive distribution

$$p(\mathcal{D}_*|\mathcal{D}) \approx p(\mathcal{D}_*|\hat{\theta})$$

Because now the posterior is just a point mass at  $\hat{\theta}$ 

In the posterior, not



PPD is more robust

(less chance of

## A "Shortcut": PPD using Marginal Likelihood

PPD, by definition, is obtained by the following marginalization

 $p(\mathcal{D}_*|\mathcal{D}) = \int p(\mathcal{D}_*|\theta) p(\theta|\mathcal{D}) d\theta$ 

Can also compute PPD without computing the posterior! Some ways:



- 2. If  $p(\mathcal{D}_*|\mathcal{D})$  can be obtained easily from the joint  $p(\mathcal{D}_*, \mathcal{D})$ 
  - Note that the PPD  $p(\mathcal{D}_*|\mathcal{D})$  is also a conditional distribution

Will see this being used we we study Gaussian Process (GP)

**CS772A: PML** 

For some distributions (e.g., Gaussian), conditionals can be easily derived from joint

# Bernoulli Observation Model



#### Estimating a Coin's Bias

- Consider a sequence of N coin toss outcomes (observations)
- Each observation  $y_n$  is a binary random variable. Head:  $y_n = 1$ , Tail:  $y_n = 0$

• Each  $y_n$  is assumed generated by a **Bernoulli distribution** with param  $\theta \in (0,1)$ Likelihood or

observation model  $p(y_n|\theta) = \text{Bernoulli}(y_n|\theta) = \theta^{y_n} (1-\theta)^{1-y_n}$ 

- Here  $\theta$  the unknown param (probability of head). Let's do MLE assuming i.i.d. data
- Log-likelihood:  $\sum_{n=1}^{N} \log p(y_n | \theta) = \sum_{n=1}^{N} [y_n \log \theta + (1 y_n) \log (1 \theta)]$
- Maximizing log-lik, or minimizing neg. log-lik (NLL) w.r.t. θ gives



I tossed a coin 5 times – gave 1 head and  
4 tails. Does it means 
$$\theta = 0.2?$$
? The  
MLE approach says so. What is I see 0  
head and 5 tails. Does it mean  $\theta = 0$ ?  
 $\theta_{MLE} = \frac{\sum_{n=1}^{N} y_n}{N}$ 

Thus MLE solution is simply the fraction of heads! <sup>(2)</sup> Makes intuitive sense! Indeed, with a small number of training observations, MLE may overfit and may not be reliable. An alternative is MAP estimation which can incorporate a prior distribution over  $\theta$ 

Probability

of a head

#### Estimating a Coin's Bias

- Let's do MAP estimation for the bias of the coin
- Each likelihood term is Bernoulli

$$p(y_n|\theta) = \text{Bernoulli}(y_n|\theta) = \theta^{y_n} (1-\theta)^{1-y_n}$$

- Also need a prior since we want to do MAP estimation
- Since  $\theta \in (0,1)$ , a reasonable choice of prior for  $\theta$  would be Beta distribution



#### Estimating a Coin's Bias

The log posterior for the coin-toss model is log-lik + log-prior

$$LP(\theta) = \sum_{n=1}^{N} \log p(y_n|\theta) + \log p(\theta|\alpha,\beta)$$

Plugging in the expressions for Bernoulli and Beta and ignoring any terms that don't depend on \u03c6, the log posterior simplifies to

$$LP(\theta) = \sum_{n=1}^{N} [y_n \log \theta + (1 - y_n) \log(1 - \theta)] + (\alpha - 1) \log \theta + (\beta - 1) \log(1 - \theta)$$

Maximizing the above log post. (or min. of its negative) w.r.t. θ gives

Using  $\alpha = 1$  and  $\beta = 1$  gives us the same solution as MLE

Recall that  $\alpha = 1$  and  $\beta = 1$  for Beta distribution is in fact equivalent to a uniform prior (hence making MAP equivalent to MLE)

 $\theta_{MAP} = \frac{\sum_{n=1}^{N} y_n + \alpha - 1}{N + \alpha + \beta - 2}$ 

Such interpretations of prior's hyperparameters as being "pseudo-observations" exist for various other prior distributions as well (in particular, distributions belonging to "exponential family" of distributions Prior's hyperparameters have an interesting interpretation. Can think of  $\alpha - 1$  and  $\beta - 1$  as the number of heads and tails, respectively, before starting the coin-toss experiment (akin to "pseudo-observations")

**CS772A: PML** 

#### The Posterior Distribution

- Let's do fully Bayesian inference and compute the posterior distribution
- Bernoulli likelihood:  $p(y_n|\theta) = \text{Bernoulli}(y_n|\theta) = \theta^{y_n} (1-\theta)^{1-y_n}$

Beta prior: 
$$p(\theta) = \text{Beta}(\theta | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha - 1} (1 - \theta)^{\beta - 1}$$
 Number of tails  $(N_0)$ 
The posterior can be computed as
$$\theta^{\sum_{n=1}^{N} y_n} (1 - \theta)^{N - \sum_{n=1}^{N} y_n} (1 - \theta)^{N - \sum_{n=1}^{N} y_n} p(\theta) p(\theta) = \frac{p(\theta) \prod_{n=1}^{N} p(y_n | \theta)}{p(y)} = \frac{p(\theta) \prod_{n=1}^{N} p(y_n | \theta)}{p(y)}$$

Here, even without computing the denominator (marg lik), we can identify the posterior

Hint: Use the fact that the

posterior must integrate to 1

 $\int p(\theta | \mathbf{y}) d\theta = 1$ 

• It is Beta distribution since  $p(\theta|\mathbf{y}) \propto \theta^{\alpha+N_1-1}(1-\theta)^{\beta+N_0-1}$ 

• Thus 
$$p(\theta|\mathbf{y}) = \text{Beta}(\theta|\alpha + N_1, \beta + N_0)$$

- Here, finding the posterior boiled down to simply "multiply, add stuff, and identify"
- Here, posterior has the same form as prior (both Beta): property of conjugate priors. PML

Exercise: Show that the

normalization constant equals

 $\Gamma(\alpha + \beta + N)$ 

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

### Conjugacy and Conjugate Priors

- Many pairs of distributions are conjugate to each other
  - Bernoulli (likelihood) + Beta (prior) ⇒ Beta posterior
  - Binomial (likelihood) + Beta (prior)  $\Rightarrow$  Beta posterior
  - Multinomial (likelihood) + Dirichlet (prior)  $\Rightarrow$  Dirichlet posterior
  - Poisson (likelihood) + Gamma (prior)  $\Rightarrow$  Gamma posterior
  - Gaussian (likelihood) + Gaussian (prior)  $\Rightarrow$  Gaussian posterior
  - and many other such pairs ..
- Tip: If two distr are conjugate to each other, their functional forms are similar
  - Example: Bernoulli and Beta have the forms

Bernoulli
$$(y|\theta) = \theta^y (1-\theta)^{1-y}$$

Beta
$$(\theta | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha - 1} (1 - \theta)^{\beta - 1}$$

This is why, when we multiply them while computing the posterior, the exponents get added and we get the same form for the posterior as the prior but with just updated hyperparameter. Also, we can identify the posterior and its hyperparameters simply by inspection

CS772A: PML

More on conjugate priors when we look at exponential family distributions

Not true in general, but in some cases (e.g., the variance of the Gaussian likelihood is fixed)

#### Predictive Distribution

- Suppose we want to compute the prob that the next outcome  $y_{N+1}$  will be head (=1)
- The posterior predictive distribution (averaging over all  $\theta$ 's weighted by their respective posterior probabilities)

## Multinoulli Observation Model



## The Posterior Distribution MLE/MAP left as

• Assume N discrete obs  $y = \{y_1, y_2, \dots, y_N\}$  with each  $y_n \in \{1, 2, \dots, K\}$ , e.g.,

- $y_n$  represents the outcome of a dice roll with K faces
- $y_n$  represents the class label of the  $n^{th}$  example in a classification problem (total K classes)
- $y_n$  represents the identity of the  $n^{th}$  word in a sequence of words
- Assume likelihood to be multinoulli with unknown params  $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_K]$  $p(y_n | \pi) = \text{multinoulli}(y_n | \pi) = \prod_{k=1}^K \pi_k^{\mathbb{I}[y_n = k]}$  Generalization of Bernoulli to K > 2 discrete outcomes
- $\pi$  is a vector of probabilities ("probability vector"), e.g.,
  - Biases of the K sides of the dice
  - Prior class probabilities in multi-class classification  $(p(y_n = k) = \pi_k)$
  - Probabilities of observing each word of the K words in a vocabulary
- Assume a conjugate prior (Dirichlet) on  $\pi$  with hyperparams  $\alpha = [\alpha_1, \alpha_2, ..., \alpha_K]$

$$p(\boldsymbol{\pi}|\boldsymbol{\alpha}) = \mathsf{Dirichlet}(\boldsymbol{\pi}|\alpha_1, \dots, \alpha_K) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1} = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional probability}}} \prod_{k=1}^K \pi_k^{\alpha_k - 1} - \underbrace{\mathsf{Generalization of Beta to}}_{\substack{K-\text{dimensional p$$

These sum to 1

Called the

concentration

parameter of the

known for now)

Dirichlet (assumed

Large values of  $\alpha$  will give a Dirichlet peaked

around its mean (next

Each  $\alpha_k \ge 0$ 

slides illustrates this)

#### Brief Detour: Dirichlet Distribution

- An important distribution. Models non-neg. vectors  $\pi$  that also sum to one
- A random draw from K-dim Dirich. will be a point under (K-1)-dim probability simplex



Basically, probability vectors

#### Brief Detour: Dirichlet Distribution

• A visualization of Dirichlet distribution for different values of concentration param



• Interesting fact: Can generate a K-dim Dirichlet random variable by independently generating K gamma random variables and normalizing them to sum to 1 CS772A: PML

#### The Posterior Distribution

• Posterior  $p(\boldsymbol{\pi}|\boldsymbol{y})$  is easy to compute due to conjugacy b/w multinoulli and Dir.

 $p(\boldsymbol{\pi}|\boldsymbol{y},\boldsymbol{\alpha}) = \frac{p(\boldsymbol{\pi},\boldsymbol{y}|\boldsymbol{\alpha})}{p(\boldsymbol{y}|\boldsymbol{\alpha})} = \frac{p(\boldsymbol{\pi}|\boldsymbol{\alpha})p(\boldsymbol{y}|\boldsymbol{\pi},\boldsymbol{\alpha})}{p(\boldsymbol{y}|\boldsymbol{\alpha})} = \frac{p(\boldsymbol{\pi}|\boldsymbol{\alpha})p(\boldsymbol{y}|\boldsymbol{\pi})}{p(\boldsymbol{y}|\boldsymbol{\alpha})} = \frac{p(\boldsymbol{\pi}|\boldsymbol{\alpha})p(\boldsymbol{y}|\boldsymbol{\pi})}{p(\boldsymbol{y}|\boldsymbol{\alpha})}$ Don't need to compute for this case because of conjugacy  $p(\boldsymbol{y}|\boldsymbol{\alpha}) = \frac{p(\boldsymbol{y}|\boldsymbol{\alpha})p(\boldsymbol{y}|\boldsymbol{\alpha})}{p(\boldsymbol{y}|\boldsymbol{\alpha})}$ Don't need to compute for this case because of conjugacy  $p(\boldsymbol{y}|\boldsymbol{\alpha}) = \frac{p(\boldsymbol{y}|\boldsymbol{\alpha})p(\boldsymbol{y}|\boldsymbol{\alpha})}{p(\boldsymbol{y}|\boldsymbol{\alpha})}$ 

Prior

estimation problem

**CS772A: PML** 

Likelihood

• Assuming  $y_n$ 's are i.i.d. given  $\boldsymbol{\pi}$ ,  $p(\boldsymbol{y}|\boldsymbol{\pi}) = \prod_{n=1}^N p(y_n|\boldsymbol{\pi})$ , and therefore  $p(\boldsymbol{\pi}|\boldsymbol{y}, \boldsymbol{\alpha}) \propto \prod_{k=1}^K \pi_k^{\alpha_k - 1} \times \prod_{n=1}^N \prod_{k=1}^K \pi_k^{\mathbb{I}[y_n = k]} = \prod_{k=1}^K \pi_k^{\alpha_k + \sum_{n=1}^N \mathbb{I}[y_n = k] - 1}$ 

- Even without computing marg-lik,  $p(y|\alpha)$ , we can see that the posterior is Dirichlet
- Denoting  $N_k = \sum_{n=1}^N \mathbb{I}[y_n = k]$ , number of observations with with value k $p(\pi|y,\alpha) = \text{Dirichlet}(\pi|\alpha_1 + N_1, \alpha_2 + N_2, \dots, \alpha_K + N_K)$ Similar to number of heads and tails for the coin bias

• Note:  $N_1$ ,  $N_2$  ...,  $N_K$  are the sufficient statistics for this estimation problem

 We only need the suff-stats to estimate the parameters and values of individual observations aren't needed (another property from exponential family of distributions – more on this later)

#### The Predictive Distribution

Thus PPD

- Finally, let's also look at the posterior predictive distribution for this model
- PPD is the prob distr of a new  $y_* \in \{1, 2, \dots, K\}$ , given training data  $y = \{y_1, y_2, \dots, y_N\}$ Will be a multinoulli. Just need  $-p(y_*|\mathbf{y}, \boldsymbol{\alpha}) = \int p(y_*|\boldsymbol{\pi}) \boldsymbol{p}(\boldsymbol{\pi}|\mathbf{y}, \boldsymbol{\alpha}) d\boldsymbol{\pi}$ to estimate the probabilities of each of the *K* outcomes
- $p(y_*|\boldsymbol{\pi}) = \text{multinoulli}(y_*|\boldsymbol{\pi}), \ p(\boldsymbol{\pi}|\boldsymbol{y},\boldsymbol{\alpha}) = \text{Dirichlet}(\boldsymbol{\pi}|\alpha_1 + N_1, \alpha_2 + N_2, \dots, \alpha_K + N_K)$
- Can compute the posterior predictive <u>probability</u> for each of the K possible outcomes

$$p(y_* = k | y, \alpha) = \int p(y_* = k | \pi) p(\pi | y, \alpha) d\pi$$
  

$$= \int \pi_k \times \text{Dirichlet}(\pi | \alpha_1 + N_1, \alpha_2 + N_2, ..., \alpha_K + N_K) d\pi$$
  

$$= \frac{\alpha_k + N_k}{\sum_{k=1}^{K} \alpha_k + N} \quad \text{(Expectation of } \pi_k \text{ w.r.t the Dirichlet posterior)}$$
  
Thus PPD is multinoulli with probability vector  $\left\{ \frac{\alpha_k + N_k}{\sum_{k=1}^{K} \alpha_k + N} \right\}_{k=1}^{K} \quad \text{Note how these probabilities have been "smoothened" due to the use of the prior + the averaging over the posterior}$   
Plug-in predictive will also be multinoulli but with prob vector given by the point estimate of  $\pi$ 

**CS772A: PML**