Probabilistic Modeling of Data: The Basics

CS771: Introduction to Machine Learning Pivush Rai

The Probabilistic Approach to ML

- Many ML problems can be seen as estimating a probability distribution/density
- Sup. Learning: Given labelled data $(X, y) = \{(x_i, y_i)\}_{i=1}^N$, estimate p(y|x)
- Unsup. Learning: Given unlabelled data $X = \{x_i\}_{i=1}^N$, estimate p(x)
- We estimate these using the given training data
 - These distributions will have some parameters θ (to be estimated)
 - These distributions will typically have a known form (which we will assume, e.g., Gaussian), but sometimes not (i.e., the form itself may also need to be estimated)
 Distribution of test data conditioned on the training data
- Once these are estimated, we can compute predictive distributions, e.g.,
 - Sup. Learning: Given a new test input x_* , what is $p(y_*|x_*, X, y)$, or mean/variance of y_* ?
 - Unsup. Learning: Given a new test input \boldsymbol{x}_* , what is $p(\boldsymbol{x}_*|\boldsymbol{X})$?



p(y=green|x)

p(y=red|x)

Getting Started: A Simple Setting

E.g., outcomes of N coin tosses, or heights of Nstudents in a class

- Assume we are given N observations $y = \{y_1, y_2, \dots, y_N\}$
- Assume these are generated from a probability model (a distribution)

 $y_n \sim p(y|\theta)$ $\forall n$ (assumed independently & identically distributed (i.i.d.))

New test observation

- Assume the form of $p(y|\theta)$ to be known (e.g., Bernoulli or Gaussian) and parameters θ of this distribution to be unknown
- Estimating θ here means we are estimating the distribution
- We can perform estimation of θ in two ways
 - Its single best/optimal value (called "point estimate")
 - A set/distribution of likely values
- Finally, we may be interested in the predictive distribution $p(y_*|y)$



Parameter Estimation in Probabilistic Models

Since data is assumed to be i.i.d., we can write down its total probability as

$$p(\mathbf{y}|\theta) = p(y_1, y_2, \dots, y_N|\theta) = \prod_{n=1}^N p(y_n|\theta)$$

• $p(y|\theta)$ called "likelihood" - probability of observed data as a function of params θ

- We wish to find the "best" heta, given observed data $m{y}$

Basically, which value of θ makes the observed data most probable under the assumed distribution $p(y|\theta)$

• One notion of "best" is to find θ which <u>maximizes</u> the likelihood



Maximum Likelihood Estimation (MLE)

- \blacksquare The goal in MLE is to find the optimal θ by maximizing the likelihood
- In practice, we maximize the log of the likelihood (log-likelihood in short)



Thus the MLE problem is

$$\theta_{MLE} = \operatorname{argmax}_{\theta} LL(\theta) = \operatorname{argmax}_{\theta} \sum_{n=1}^{N} \log p(y_n | \theta)$$

This is now an optimization (maximization problem)



CS771: Intro to ML

• The negative log-lik $(-\log p(y_n|\theta))$ is akin to the loss on each data point

Such priors have various other benefits as we will see later

Negative Log-Likelihood

(NLL)



- $\theta_{MLE} = \operatorname{argmax}_{\theta} \sum_{n=1}^{N} \log p(y_n | \theta) = \operatorname{argmin}_{\theta} \sum_{n=1}^{N} -\log p(y_n | \theta)$
- Thus MLE can also be seen as minimizing the negative log-likelihood (NLL)

Maximum Likelihood Estimation (MLE)

The MLE problem can also be easily written as a <u>minimization</u> problem

$$\theta_{MLE} = \operatorname*{argmin}_{\theta} NLL(\theta)$$

NLL is analogous to a loss function

Thus doing MLE is akin to <u>minimizing training loss</u>





MLE: An Example

- 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 $\theta \in (0,1)$

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

- Here θ the unknown param (probability of head). Want to estimate it using MLE
- 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 NLL) w.r.t. θ will give a closed form expression



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! \textcircled{O} Makes
intuitive sense!
Indeed – if you want to trust
MLE solution. But with small
number of training
observations, MLE may overfin
and may not be reliable. We want to trust
that use prior distributions!



Intro to ML

Probability

of a head

Take deriv. set it

to zero and solve. Easy optimization

MLE and Its Shortcomings..

MLE finds parameter values s that make the observed data most probable

$$\theta_{MLE} = \arg\max_{\theta} \sum_{n=1}^{N} \log p(y_n | \theta) = \arg\min_{\theta} \sum_{n=1}^{N} -\log p(y_n | \theta)$$

- No provision to control overfitting (MLE is just like minimizing training loss)
- How do we regularize probabilistic models in a principled way?
- Also, MLE gives only a single "best" answer ("point estimate")
 - .. and it may not be very reliable, especially when we have very little data
 - Desirable: Report a probability distribution over the learned params instead of point est
- Prior distributions provide a nice way to accomplish such things!

This distribution can give us a sense about the <u>uncertainty</u> in the parameter estimate



8

Before observing any data

Can specify our prior belief about likely param values via a prob. dist., e.g.,

This is a rather simplistic/contrived prior. OJust to illustrate the basic idea. We will see more concrete examples of priors shortly. Also, the prior usually depends (assumed conditioned on) on some fixed/learnable hyperparameters (say some α and β , and written as $p(\theta | \alpha, \beta)$

Priors





Fully Bayesian

inference

Once we observe the data y, apply Bayes rule to update prior into posterior

 $\frac{Prior}{p(\theta)p(\mathbf{y}|\theta)}$

Likelihood
 Note: Marginal lik. is hard to compute in general as it requires a summation or integral which may not be easy (will briefly look at this in CS771, although will stay away going too deep in this course – CS772 does that in more detail)

Maximum-a-

Two ways now to report the answer:

Posterior

• Report the maxima (mode) of the posterior: $\arg \max_{\theta} p(\theta | \mathbf{y}) \prec \det_{\text{estimation}}^{\text{posteriori} (MAP)}$

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

Report the full posterior (and its properties, e.g., mean, mode, variance, quantiles set () ntro to ML

Posterior

- Posterior distribution tells us how probable different parameter values are <u>after</u> we have observed some data
- Height of posterior at each value gives the posterior probability of that value



- \star More likely values
- \star Less likely values

Can think of the posterior as a "hybrid" obtained by combining information from the likelihood and the prior

CS771: Intro to ML

Maximum-a-Posteriori (MAP) Estimation

The MAP estimation approach reports the maxima/mode of the posterior

$$\theta_{MAP} = \arg\max_{\theta} p(\theta|y) = \arg\max_{\theta} \log p(\theta|y) = \arg\max_{\theta} \log \frac{p(\theta)p(y|\theta)}{p(y)}$$

• Since p(y) is constant w.r.t. θ , the above simplifies to

$$\theta_{MAP} = \arg \max_{\theta} \left[\log p(y|\theta) + \log p(\theta) \right]$$

$$= \arg\min_{\theta} \left[-\log p(y|\theta) - \log p(\theta) \right]$$

The NLL term acts like the training loss and the (negative) log-prior acts as regularizer. Keep in mind this analogy. ©

 $\theta_{MAP} = \arg \min_{\theta} \left[NLL(\theta) - \log p(\theta) \right]$

- Same as MLE with an extra log-prior-distribution term (acts as a regularizer) ☺
- If the prior is absent or <u>uniform</u> (all values equally likely a prior) then MAP=MLE



CS771: Intro to ML

MAP Estimation: An Example

- Let's again consider the coin-toss problem (estimating 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 θ would be Beta distribution



MAP Estimation: An Example (Contd)

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

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

 $\hfill Plugging in the expressions for Bernoulli and Beta and ignoring any terms that don't depend on <math display="inline">\theta$, 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")

CS771: Intro to ML

Fully Bayesian Inference

MLE/MAP only give us a point estimate of

MAP estimate is more robust than MLE (due to the regularization effect) but the estimate of uncertainty is missing in both approaches – both just return a single "optimal" solution by solving an optimization problem



Interesting fact to keep in mind: Note that the use of the prior is making the MLE solution move towards the prior (MAP solution is kind of a "compromise between MLE solution of the mode of the prior) ©



MI

14

Fully Bayesian inference

If we want more than just a point estimate, we can compute the full posterior

Computable analytically only when the prior and likelihood are "friends" with each other (i.e., they form a conjugate pair of distributions (distributions from exponential family have conjugate priors

$$p(\theta|\mathbf{y}) = \frac{p(\theta)p(\mathbf{y}|\theta)}{p(\mathbf{y})}$$

An example: Bernoulli and Beta are conjugate. Will see some more such pairs

In other cases, the posterior needs to be approximated (will see 1-2 such cases in this course; more detailed treatment in the advanced course on probabilistic modeling and inference)

"Online" Nature of Bayesian Inference

Fully Bayesian inference fits naturally into an "online" learning setting



Also, the posterior becomes more and more "concentrated" as the number of observations increases. For very large N, you may expect it to be peak around the MLE solution



• Our belief about θ keeps getting updated as we see more and more data



15



Conjugacy

- 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



Probabilistic Models: Making Predictions

- Having estimated θ , we can now use it to make predictions
- Prediction entails computing the predictive distribution of a new observation, say y_*

$$p(y_{*}|y) = \int p(y_{*},\theta|y) d\theta \qquad \text{Marginalizing over the unknown } \theta$$
Conditional distribution of the new observation, given past observations
$$= \int p(y_{*}|\theta,y)p(\theta|y) d\theta \qquad \text{Decomposing the joint using chain rule}$$

$$= \int p(y_{*}|\theta)p(\theta|y) d\theta \qquad \text{Assuming i.i.d. data, given } \theta, y_{*} \text{ does not depend on } y$$

- When doing MLE/MAP, we approximate the posterior $p(\theta|\mathbf{y})$ by a single point θ_{opt} $p(y_*|\mathbf{y}) = \int p(y_*|\theta)p(\theta|\mathbf{y})d\theta \approx p(y_*|\theta_{opt})$ A "plug-in prediction" (simply plugged in the single best estimate we had)
- When doing fully Bayesian estimation, getting the predictive dist. will require computing

$$p(y_*|\mathbf{y}) = \int p(y_*|\theta) p(\theta|\mathbf{y}) d\theta$$
$$\mathbb{E}_{p(\theta|\mathbf{y})}[p(y_*|\theta)]$$

This computes the predictive distribution by averaging over the full posterior – basically calculate $p(y_*|\theta)$ for each possible θ , weighs it by how likely this θ is under the posterior $p(\theta|y)$, and sum all such posterior weighted predictions. Note that not each value of theta is given equal importance here in the averaging



ntro to ML

For example, PMF of the label of a

new test input in classification

Probabilistic Models: Making Predictions (Example)

- For coin-toss example, let's compute probability of the $(N + 1)^{th}$ toss showing head
- This can be done using the MLE/MAP estimate, or using the full posterior

$$\theta_{MLE} = \frac{N_1}{N}$$
 $\theta_{MAP} = \frac{N_1 + \alpha - 1}{N + \alpha + \beta - 2}$ $p(\theta|\mathbf{y}) = \text{Beta}(\theta|\alpha + N_1, \beta + N_0)$

Thus for this example (where observations are assumed to come from a Bernoulli)

MLE prediction: $p(y_{N+1} = 1|\mathbf{y}) = \int p(y_{N+1} = 1|\theta)p(\theta|\mathbf{y})d\theta \approx p(y_{N+1} = 1|\theta_{MLE}) = \theta_{MLE} = \frac{N_1}{N}$ MAP prediction: $p(y_{N+1} = 1|\mathbf{y}) = \int p(y_{N+1} = 1|\theta)p(\theta|\mathbf{y})d\theta \approx p(y_{N+1} = 1|\theta_{MAP}) = \theta_{MAP} = \frac{N_1 + \alpha - 1}{N + \alpha + \beta - 2}$

Fully Bayesian:
$$p(y_{N+1} = 1|\mathbf{y}) = \int p(y_{N+1} = 1|\theta)p(\theta|\mathbf{y})d\theta = \int \theta p(\theta|\mathbf{y})d\theta = \int \theta Beta(\theta|\alpha + N_1, \beta + N_0)d\theta = \frac{N_1 + \alpha}{N + \alpha + \beta}$$

Expectation of θ under the Beta

Bayesian inference

posterior that we computed using fully

Sector to ML



Again, keep in mind that the posterior weighted averaged prediction used in the fully Bayesian case would usually not be as simple to compute as it was in this case. We will look at some hard cases later

Probabilistic Modeling: A Summary

- Likelihood corresponds to a loss function; prior corresponds to a regularizer
- Can choose likelihoods and priors based on the nature/property of data/parameters
- MLE estimation = unregularized loss function minimization
- MAP estimation = regularized loss function minimization
- Allows us to do fully Bayesian learning (learning the full distribution of the parameters)
- Makes robust predictions by posterior averaging (rather than using point estimate)
- Many other benefits, such as
 - Estimate of confidence in the model's prediction (useful for doing Active Learning)
 - Can do automatic model selection, hyperparameter estimation, handle missing data, etc.
 - Formulate latent variable models
 - .. and many other benefits (a proper treatment deserves a separate course, but we will see some of these in this course, too)

CS771: Intro to ML