Bayesian Machine Learning for Scientific Research

Maximizing information from data

Ezequiel Alvarez sequi@unsam.edu.ar ICTP-SAIFR October 2024

Previous lectures

Lecture 1



Bayesian: assume data being sampled from a PDF, infer its parameters and learn the internal structure of the data

> To learn the PDF of the data and then being able to assess, predict , generate, etc.

> > More scientific



Bayes Theorem:

 $p(\theta \mid X) = \underline{p(X \mid \theta) * p(\theta)}$ p(x)

Our utility: X = data, θ=parameters

Model data as being sampled from a clever PDF with parameters θ

Infer θ once you see the data X

Connect θ to physical parameters of interest



1D Gaussian Mixture

Lecture 2: Mixture of Bernoulli



Input



Scientifically: z_n is the probability of each class

Lecture 3: Mixture Models



$$p(x_n|\Theta) = \sum_{k=1}^K \pi_k \, p(x_n| heta_k)$$
 $ext{g} p(X|\Theta) = \sum_{n=1}^N \log(\sum_{k=1}^K \pi_k \, p(x| heta_k) \,)$

Mixture Models: Crypto price ?



Real Bitcoin price Aug to Sep

sigmas = [10,100,1000]
<pre>crypto_price = [50000]</pre>
np.random.seed(1)
<pre>z = np.random.choice([0,1,2], p=[0.7,0.25,0.05], size=5000)</pre>
for zn in z:
<pre>crypto price.append(cripto price[-1]+np.random.normal(0,sigmas[zn]))</pre>

plt.plot(cripto_price)
plt.title('Emulated crypto price')
plt.ylabel('USD')
plt.xlabel('Time')
plt.show()



Lecture 3: Dirichlet Distribution



Lecture 3: Mixture Models, where is the hack?





Lecture 3: 2D Mixture Model, hh \rightarrow **bbyy**



Lecture 3: The impossible....@10%





Lecture 3: The impossible....@5%!

0.00





Assessment in Bayesian ML Lecture 4

Assessment in Bayesian ML



Once you have a result from real data, you want to test that...

- The sampling is unbiased
- The model works correctly
- The model is (fairly) correct for the data

Assessment in Bayesian ML



Once you have a result from real data, you want to test that...

• The sampling is unbiased

Sampling diagnostics

• The model works correctly

Fake data to test it

• The model is (fairly) correct for the data

Posterior Predictive Check

Sampling Assessment

@Hamiltonian MonteCarlo

MCMC Sampling with Hamiltonian MC $p(\theta)$ $V(\theta) = -ln(p(\theta))$



 $V(\theta) = -\ln(p(\theta))$ p(θ)

We know all its gradients



Simulate trajectory of fictitious particle



Simulate trajectory of fictitious particle

Introduce stochastic energy loss in order to stay more time in the larger probability regions



parameters	lp	accept_stat	stepsize	treedepth	n_leapfrog	divergent	energy	mu	sigma
draws									
0	-32598.734453	0.994393	0.176560	3.0	11.0	0.0	32600.239461	120.773958	6.188426
1	-32601.510781	0.808358	0.149106	4.0	15.0	0.0	32604.731098	115.727014	10.256436
2	-32597.717577	0.973526	0.137055	3.0	7.0	0.0	32599.468725	120.248254	6.372511
3	-32597.238541	0.998011	0.141012	4.0	27.0	0.0	32599.385174	119.383923	7.435660
4	-32598.284739	0.988675	0.176560	4.0	15.0	0.0	32600.078479	119.809050	6.475248

p(sample) _____ HMC numbers_____/ parameters

lp	accept_stat	stepsize	treedepth	n_leapfrog	divergent	energy	mu	sigma
-32598.734453	0.994393	0.176560	3.0	11.0	0.0	32600.239461	120.773958	6.188426
-32601.510781	0.808358	0.149106	4.0	15.0	0.0	32604.731098	115.727014	10.256436
-32597.717577	0.973526	0.137055	3.0	7.0	0.0	32599.468725	120.248254	6.372511
-32597.238541	0.998011	0.141012	4.0	27.0	0.0	32599.385174	119.383923	7.435660
-32598.284739	0.988675	0.176560	4.0	15.0	0.0	32600.078479	119.809050	6.475248
	lp -32598.734453 -32601.510781 -32597.717577 -32597.238541 -32598.284739	lp_ accept_stat_ -32598.734453 0.994393 -32601.510781 0.808358 -32597.717577 0.973526 -32597.238541 0.998011 -32598.284739 0.988675	lp_ accept_stat_ stepsize_ -32598.734453 0.994393 0.176560 -32601.510781 0.808358 0.149106 -32597.717577 0.973526 0.137055 -32597.238541 0.998011 0.141012 -32598.284739 0.988675 0.176560	lp_ accept_stat_ stepsize_ treedepth_ -32598.734453 0.994393 0.176560 3.0 -32601.510781 0.808358 0.149106 4.0 -32597.717577 0.973526 0.137055 3.0 -32597.238541 0.998011 0.141012 4.0 -32598.284739 0.988675 0.176560 4.0	lp_ accept_stat_ stepsize_ treedepth_ n_leapfrog_ -32598.734453 0.994393 0.176560 3.0 11.0 -32601.510781 0.808358 0.149106 4.0 15.0 -32597.717577 0.973526 0.137055 3.0 7.0 -32597.238541 0.998011 0.141012 4.0 27.0 -32598.284739 0.988675 0.176560 4.0 15.0	lp_ accept_stat_ stepsize_ treedepth_ n_leapfrog_ divergent_ -32598.734453 0.994393 0.176560 3.0 11.0 0.0 -32601.510781 0.808358 0.149106 4.0 15.0 0.0 -32597.717577 0.973526 0.137055 3.0 7.0 0.0 -32597.238541 0.998011 0.141012 4.0 27.0 0.0 -32598.284739 0.988675 0.176560 4.0 15.0 0.0	lpaccept_stat_stepsizetreedepthn_leapfrogdivergent_energy32598.7344530.9943930.1765603.011.00.032600.239461-32601.5107810.8083580.1491064.015.00.032604.731098-32597.7175770.9735260.1370553.07.00.032599.468725-32597.2385410.9980110.1410124.027.00.032599.385174-32598.2847390.9886750.1765604.015.00.032600.078479	lpaccept_stat_stepsizetreedepth_n_leapfrogdivergent_energy_mu-32598.7344530.9943930.1765603.011.00.032600.239461120.773958-32501.5107810.8083580.1491064.015.00.032604.731098115.727014-32597.7175770.9735260.1370553.07.00.032599.468725120.248254-32597.2385410.9980110.1410124.027.00.032599.385174119.383923-32598.2847390.9886750.1765604.015.00.032600.078479119.809050

p(sample) _____ HMC numbers_____/ parameters

Simple solutions, e.g. make a more detailed exploration

fit = model.sample(num_chains=4, num_samples=1000, delta=0.95)



Check that all (independent) chains end up exploring approximately the same!



Check that all (independent) chains end up exploring approximately the same!

R-hat < 1.05

R-hat = sqrt((Variance Between Chains) / (Variance Within Chains))

<pre>summary = az.summary(azdata) summary</pre>										
	mean	sd	hdi_3%	hdi_97%	mcse_mean	mcse_sd	ess_bulk	ess_tail	r_hat	
mu	118.861	1.320	116.425	120.933	0.055	0.039	691.0	570.0	1.01	
sigma	7.770	1.120	5.868	9.878	0.046	0.033	660.0	702.0	1.01	
lambda0	0.062	0.004	0.054	0.069	0.000	0.000	1114.0	842 <mark>.</mark> 0	1.00	
theta[0]	0.348	0.089	0.200	0.505	0.004	0.003	630.0	678.0	1.01	
theta[1]	0.652	0.089	0.495	0.800	0.004	0.003	630.0	678.0	1.01	

azdata - az from nystan(fit)

Check that all (independent) chains end up exploring approximately the same!

R-hat < 1.05

R-hat = sqrt((Variance Between Chains) / (Variance Within Chains))

<pre>summary = az.summary(azdata) summary</pre>										
	mean	sd	hdi_3%	hdi_97%	mcse_mean	mcse_sd	ess_bulk	ess_tail	r_hat	
mu	118.861	1.320	116.425	120.933	0.055	0.039	691.0	570.0	1.01	
sigma	7.770	1.120	5.868	9.878	0.046	0.033	660.0	702.0	1.01	
lambda0	0.062	0.004	0.054	0.069	0.000	0.000	1114.0	<mark>842.</mark> 0	1.00	
theta[0]	0.348	0.089	0.200	0.505	0.004	0.003	630.0	678.0	1.01	
theta[1]	<mark>0.652</mark>	0.089	0.495	0.800	0.004	0.003	630.0	678.0	1.01	

azdata = az from nystan(fit)

Check that all (independent) chains end up exploring approximately the same!

R-hat < 1.05

R-hat = sqrt((Variance Between Chains) / (Variance Within Chains))

There are a few more diagnostic indicators, but this plus some adjustments in adapt, max_tree, etc. is usually enough

Never use a model before testing it

Never use a model before testing it

• Create fake data as you expect it to be the real data

Never use a model before testing it

- Create fake data as you expect it to be the real data
- You know the trues, start then with biased priors

Never use a model before testing it

- Create fake data as you expect it to be the real data
- You know the trues, start then with biased priors
- Verify that the inference process approaches the posterior to the (known)

trues



Never use a model before testing it

- Create fake data as you expect it to be the real data
- You know the trues, start the
- Verify that the inference pro

trues











You have (at least) 2 metrics to test

- Density at the *true*
- Distance *mean* to *true*
Model works correctly?



You have (at least) 2 metrics to test

- Density at the *true*
- Distance *mean* to *true*



Posterior Predictive Check

The probability of the data

How good is the model to explain the data

Posterior predictive check

How good is the model to explain the data?

Posterior predictive check

How good is the model to explain the data? Once you infer:

What is the probability of the data within the model?

Posterior predictive check

How good is the model to explain the data? Once you infer:

What is the probability of the data within the model ?

How we measure what is good and what is bad ??

- Data is usually multidimensional

- Data is usually multidimensional
- After inference we can easily sample tons of synthetic data!

- Data is usually multidimensional
- After inference we can easily sample tons of synthetic data!



- Data is usually multidimensional
- After inference we can easily sample tons of synthetic data!



1D case

You have some data

Solve the inference problem with your model

- Data is usually multidimensional
- After inference we can easily sample tons of synthetic data!



- Data is usually multidimensional
- After inference we can easily sample tons of synthetic data!



- Data is usually multidimensional
- After inference we can easily sample tons of synthetic data!





$$p = 10^{-6} \dots and now?$$



- What does it mean ?
- What do we compare it to?



- What does it mean?
- What do we compare it to?
- Generate replicas of data X^{Rep}



- What does it mean?
- What do we compare it to?
- Generate replicas of data X^{Rep}
- Compute their probability



- What does it mean ?
- What do we compare it to?
- Generate replicas of data X^{Rep}
- Compute their probability
- Compute

 $p(p(X^{rep}) < p(X))$



• S ~ 1: bad model

Score =
$$p(p(X^{rep}) < p(X))$$

- S ~ 1 : bad model
- S ~ 0.5: good model



- S ~ 1: bad model
- S ~ 0.5: good model
- $S \lesssim 0.1$: bad model



Avoid overfit: given the data X, leave a held-out X_{held} for the posterior predictive check!

Score =
$$p(p(X_{held}^{rep}) < p(X_{held}))$$

Realistic 2D example

 $\mu = 120, \sigma = 7, \lambda = 0.06 \text{ and } \Theta_S = 0.3$



 $\mu = 120, \sigma = 7, \lambda = 0.06 \text{ and } \Theta_s = 0.3$ μ =120, σ =7, λ =0.06 and Θ_{S} =0.3 data nvariant mass [GeV] Invariant mass [GeV] Invariant mass [GeV] Invariant mass [GeV] Data

Labelled Data

Background Signal





Data: 80% to infer, 20% held-out

Held-out data

Posterior Predictive Check: pp \rightarrow bbyy

 $\mu = 120, \sigma = 7, \lambda = 0.06 \text{ and } \Theta_S = 0.3$



my model = """ data { int <lower=0> N: // number of datapoints real L, U; // lower and upper limits of the observables (i.e. how to truncate the distributions) array[N,2] real<lower=L, upper=U> y; // there are 2 observables per datapoint real mu0, sigma0, mu1, sigma1, mu2, sigma2, t1, t2; // hyperparameters for the parameters priors parameters { real<lower=0> mu: real<lower=0> sigma: real<lower=0> lambda0; simplex[2] theta; model { vector[2] lp; mu ~ normal(mu0, sigma0); sigma ~ cauchy(mul, sigmal); lambda0 ~ normal(mu2, sigma2); theta ~ dirichlet([t1,t2]); for (n in 1:N) { lp[1] = exponential lpdf(y[n,1] - L | lambda0) - exponential lcdf(U - L | lambda0) + exponential lpdf(y[n,2] - L| lambda0) - exponential lcdf(U - L | lambda0); lp[2] = normal lpdf(v[n,1] | mu, sigma)- log diff exp(normal lcdf(U | mu, sigma), normal lcdf(L | mu, sigma)) + normal lpdf(y[n,2] | mu, sigma) - log diff exp(normal lcdf(U | mu, sigma), normal lcdf(L | mu, sigma)); // Pay attention here: theta[1] is associated to background, and theta[2] to signal, since we don't have theta[0] target += log mix(theta, lp); };

We infer with our model



Posterior x model



Posterior x model = sinthetic data !







0.014
0.012
$$p(x_n|X) = \int p(x_n|z_i) \, p(z_i|X) \, dz_i$$

0.000
0.006
0.004
0.002



$$p(x_n|X) = \int p(x_n|z_i) p(z_i|X) dz_i$$
We've never seen
a label!





$$p(x_{n}|X) = \int p(x_{n}|z_{i}) p(z_{i}|X) dz_{i}$$

$$p(x_{n}|z_{i}) p(z_{i}|X) dz_{i}$$
























Probability of X_{held} to have been sampled from the inferred model!



 $p\left(p(X_{held}^{rep}) < p(X_{held})
ight) = 0.86$

Probability of X_{held} to have been 0.025 -Replicas of the data Data sampled from the inferred model! 0.020 No general way: 0.015 -Safeguard against gross 0.010 misspecification 0.005 0.000

-9520

-9500

Probability of the probability

 $p\left(p(X_{held}^{rep}) < p(X_{held})
ight) = 0.86$

-9480

Log(p)

-9460

-9440



 $p\left(p^2(X_{held}^{rep}) < p^2(X_{held})
ight) = 0.37$

Final remarks @ Lecture 4

Assessment in Bayesian Machine Learning

Final remarks

• There are ways to check an inference process

Final remarks

• There are ways to check an inference process

- Sampling diagnostics
- Constructed data to find and understand problems
- Posterior predictive checks

Final remarks

• There are ways to check an inference process

- Sampling diagnostics
- Constructed data to find and understand problems
- Posterior predictive checks

• Bayesian Workflow: 2011.01808