

6.790 Homework 1

Sept 10, 2024

Questions 1–3 are relatively stand-alone warm-ups. Questions 4–6 are more extended practice and illustrations of the ideas of this material. Question 7 requires coding. *Do not submit your code!*

There are some rhetorical questions in blue boxes. You don't need to answer them—they're just for thinking about.

Please hand in your work via Gradescope via the link at <https://gradml.mit.edu/info/homeworks/>. If you were not added to the course automatically, please use Entry Code R7RGGX to add yourself to Gradescope.

1. Latex is not required, but if you are hand-writing your solutions, please write clearly and carefully. You should include enough work to show how you derived your answers, but you don't have to give careful proofs.
2. Homework is due on Tuesday September 17 at 11PM.
3. Lateness and extension policies are described at https://gradml.mit.edu/info/class_policy/.

Contents

1 Normal fish [10 Points]	3
1.1 Fish tale	3
1.2 A simple model	3
1.3 A more complex model	4
2 Parameter estimation [10 points]	4
2.1 Force field	4
2.2 Pigeons	4
3 Bayesian belief update [10 Points]	4
3.1 Beta-Binomial practice	4
3.2 What's new?	5
4 Which dice factory? [15 points]	6
5 Emergency Room [15 Points]	6
6 Abby Normal [15 Points]	8

7 Coding Question: Two Gaussians [25 Points]	8
8 Optional	9
8.1 Throwing rocks: asymmetric loss	9
8.2 Copy that: discrete Bayes update and decision theory	9
8.3 Dirichlet Priors	10
8.4 More fun with Dirichlet	11

1 Normal fish [10 Points]

1.1 Fish tale

(Bishop 1.11) We find ourselves with a data set consisting of the measured weights of a bunch of fish caught during an afternoon of fishing. We decide to model the distribution of these weights using a Gaussian distribution.

Why might this not be a great modeling choice?

Our goal is to select parameters μ, σ^2 of the Gaussian distribution in order to maximize the likelihood of our data, $\mathcal{D} = \{x^{(1)}, \dots, x^{(n)}\}$. The parameters that maximize the log likelihood of the data, will also maximize the likelihood (due to its monotonicity) and the form is easier to deal with. Recall that the pdf of a Gaussian distribution is given by

$$p_X(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} .$$

If we assume that the process whereby we caught the fish made their weights independent and identically distributed, then

$$p(\mathcal{D} | \mu, \sigma^2) = \prod_i p_X(x^{(i)} | \mu, \sigma^2) .$$

The log likelihood function is then

$$\log p(\mathcal{D} | \mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (x^{(i)} - \mu)^2 - \frac{N}{2} \log \sigma^2 - \frac{N}{2} \log(2\pi) .$$

By setting its derivatives with respect to μ and σ^2 equal to zero and solving, verify that the maximum likelihood estimates of μ and σ are given by

$$\mu_{\text{ml}} = \frac{1}{N} \sum_{n=1}^N x^{(n)}$$

$$\sigma_{\text{ml}}^2 = \frac{1}{N} \sum_{n=1}^N (x^{(n)} - \mu_{\text{ml}})^2$$

Under what assumptions about the log likelihood function is this a valid approach for finding a global maximum?

This solution may be different than the estimator you have previously seen for σ^2 . See the discussion at the bottom of Bishop page 27 for an explanation.

1.2 A simple model

As it happens, we caught 6 mega-guppies (a tasty type of fish), with these weights:

$$\mathcal{D}_0 = \{0.9, 1, 1.1, 1.2, 3, 3.1\} .$$

We looked in the USDA handbook which told us that the variance of the weight of North American mega-guppies is $\sigma^2 = 0.5^2 = 0.25$.

Find the maximum likelihood value of μ_{ml} for \mathcal{D}_0 under this assumption. What is the data likelihood $p(\mathcal{D}_0|\mu_{\text{ml}})$?

1.3 A more complex model

Now, what if we ignore the USDA value of σ^2 and decide to estimate it ourselves? Find the maximum likelihood estimates μ_{ml} and σ_{ml}^2 of μ and σ^2 for our data set \mathcal{D}_0 . What is the data likelihood $p(\mathcal{D}_0|\mu_{\text{ml}}, \sigma_{\text{ml}}^2)$?

What are the advantages and disadvantages of this model versus the one with the fixed variance?

2 Parameter estimation [10 points]

2.1 Force field

A supervillain has our hero trapped in an invisible one-dimensional force-field (hero can only move in one dimension) and we know that the field has finite extent. Using a drone flying overhead, we make several measurements of the hero's position.

We wish to estimate the boundaries of the force-field given samples of the hero's position.

If we knew that our data are drawn uniformly from a finite interval, $[a, b]$, then we might want to find $a_{\text{ml}}, b_{\text{ml}}$ to maximize the likelihood of \mathcal{D} .

For our data set $\mathcal{D} = (x^{(1)}, x^{(2)}, \dots, x^{(n)})$, what are the maximum likelihood parameter estimates a_{ml} and b_{ml} ? What is the data likelihood $p(\mathcal{D}|a_{\text{ml}}, b_{\text{ml}})$?

Is this model of the hero data a good one? Why or why not?

2.2 Pigeons

Pigeons¹, when put in a situation where $\Pr(y = 1) = p$ and $\Pr(y = 0) = 1 - p$, will select option 1 with probability p and option 0 with probability $1 - p$. What is the expected 0-1 loss for the pigeons' decision rule? What is the optimal decision rule and its expected loss?

Actually, people² do this too!

3 Bayesian belief update [10 Points]

3.1 Beta-Binomial practice

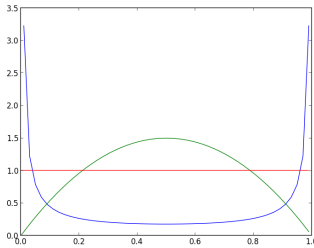
(a) Label which of the lines in the figure below correspond to:

1. Beta(0.1, 0.1)

¹"Probability-Matching in the Pigeon", Donald H. Bullock and M. E. Bitterman, *The American Journal of Psychology*, Vol. 75, No. 4 (Dec., 1962), pp. 634-639

²"Banking on a Bad Bet: Probability Matching in Risky Choice is Linked to Expectation Generation," *Psychological Science*, Vol. 22, No. 6 (2011).

2. Beta(1,1)
3. Beta(2,2)



We are estimating the probability that a coin comes up heads.

- (b) What does it mean to have a prior of Beta(2,2)?
- (c) If that's the prior, what is the posterior after seeing 3 heads and 2 tails?
- (d) What are the mean and mode of that posterior?
- (e) What does it mean to have a prior of Beta(2,3)?
- (f) If that's the prior, what is the posterior after seeing 3 heads and 2 tails?
- (g) What are the mean and mode of that posterior?

3.2 What's new?

(Bishop 2.7) Consider a bernoulli random variable x with mean μ with prior distribution for μ given by the beta distribution:

$$\text{Beta}(\mu; a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad (2.13)$$

and suppose we have observed m occurrences of $x = 1$ and l occurrences of $x = 0$. Show that the posterior mean value of μ lies between the prior mean and the maximum likelihood estimate for μ .

To do this, show that the posterior mean can be written as λ times the prior mean plus $(1 - \lambda)$ times the maximum likelihood estimate, where $0 \leq \lambda \leq 1$. This illustrates the concept of the posterior distribution being a compromise between the prior distribution and the maximum likelihood solution.

4 Which dice factory? [15 points]

You have just purchased a two-sided die, which can come up either 1 or 2:



You want to use your crazy die in some betting games with friends later this evening, but first you want to know the probability that it will roll a 1.

You know it came either from factory 0 or factory 1, but not which.

Factory 0 produces dice that roll a 1 with probability ϕ_0 . Factory 1 produces dice that roll a 1 with probability ϕ_1 . You believe initially that with probability η_0 that it came from factory 1.

- (a) Without seeing any rolls of this die, what would be your predicted probability that it would roll at 1?
- (b) If we roll the die and observe the outcome, what can we infer about where the coin was manufactured?
- (c) More concretely, let's assume that:
 - $\phi_0 = 1$: dice from factor 0 always roll a 1
 - $\phi_1 = 0.5$: dice from factory 1 are fair (roll at 1 with probability 0.5)
 - $\eta_0 = 0.7$: we think with probability 0.7 that this die came from factory 1

Now we roll it, and it comes up 1! What is your posterior distribution on which factory it came from? What is your predictive distribution on the value of the next roll?

- (d) You roll it again, and it comes up 1 again.

Now, what is your posterior distribution on which factory it came from? What is your predictive distribution on the value of the next roll?

- (e) Instead, what if it rolls a 2 on the second roll?
- (f) In the general case (not using the numerical values we have been using) prove that if you have two observations, and you use them to update your prior in two steps (first conditioning on one observation and then conditioning on the second), that no matter which order you do the updates in you will get the same result.

5 Emergency Room [15 Points]

You are a young doctor, working off your federal medical school tuition grant in southern North Dakota. It's your fourth day on the job. You are all alone in the emergency room (ER) when Pat comes in complaining of chest pain.

You have to predict whether Pat is having a heart attack (H) or indigestion (I). Your loss function is:

$$L(g, a) = \begin{cases} 0 & \text{if } g = a \\ 1 & \text{if } g = \text{"H"} \text{ and } a = \text{"I"} \\ 10 & \text{if } g = \text{"I"} \text{ and } a = \text{"H"} \end{cases}$$

You have seen three previous patients who exhibited chest pain, none of whom were actually having a heart attack.

- (a) You use those three data points to make a point estimate of the probability that Pat is having a heart attack and then use it to make the prediction that minimizes the empirical risk. What do you predict? What is the empirical risk of that prediction?

Do you think the empirical risk of this predictor is a good measure of how useful it will be?

- (b) The next morning, you think more carefully and decide it would be better to forget all your previous experience and simply view each new patient with an open mind. So, you use some ideas from this week's lectures. Let Q be a random variable representing the probability that a random patient walking into your ER will be having a heart attack. You have a uniform prior on Q .

What is the prediction that minimizes risk for a random patient walking into your ER? What is the risk of that prediction?

- (c) Later that afternoon, you figure it would be better to combine approaches. So, what if you started with a uniform prior, but then observed three patients all of whom had indigestion?

What would be your posterior distribution on Q ? What prediction should you make? What is the risk (under the posterior distribution) of that prediction?

- (d) That evening, really worried that you haven't had enough experience in these matters, and beginning to question your judgment about accepting this job, you decide to call your friend Chris who is working at Mass General. Chris has seen 20 patients with indigestion and 1 with heart attack. You use Chris's experience to construct a prior distribution, and then update it with your own (3 patients with indigestion).

What would be your posterior distribution on Q ? What prediction should you make? What is the risk (under the posterior distribution) of that prediction?

- (e) At 2AM, questioning the meaning of life, you are quite sure that you should have become a poet. You are so uncertain of your ability to make predictions that you call your former professor who is the head of the emergency medicine department at Gotham City Hospital. Herr Prof. Dr. Strangelove has seen 2000 patients with indigestion and 20 with heart attack. You use Dr. Strangelove's experience to construct a prior distribution, and then update it with your own (3 patients with indigestion).

What would be your posterior distribution on Q ? What prediction should you make? What is the risk of that prediction?

Is there a potential problem with using Dr. Strangelove's data to help construct your prior?

6 Abby Normal [15 Points]

Dr. Frahnkensteen is designing an artificial cranium, but she needs to know how big to make it; her design goal is to be a good fit to 80% of brains. So, she wants to get a good estimate of the distribution of the sizes of brains in the local population. Since brains are kind of squishy, we will just consider the total volume of the brain, a one-dimensional quantity.

The Dr. has considerable previous experience with brains and thinks their distribution is well modeled as a Gaussian distribution with with a variance of 75cc. But she's not at all sure about the mean of this current population. She thinks it might be somewhere around 1100cc.

- (a) One way to express the Dr.'s uncertainty about the distribution of brain sizes in the local population is to put a Gaussian distribution *on the mean* of the local distribution.

What are the hyper-parameters of this distribution? Pick some to model Dr. F's situation (they're not completely determined by the story).

- (b) Dr. F. sends her assistant Eygor out to get a new brain from the local population. Eygor brings back one that is 1500cc! What should the posterior be?

Start by solving this problem algebraically. Write down the prior and the observation likelihood function symbolically. Then, derive a form for the posterior.

What actual numerical values do you get, given your answer to the previous question, and the observation of 1500cc?

- (c) How is the new mean related to the old mean and the observation?
- (d) What can we say about how the variance behaves when an observation is made?
- (e) What is Dr. F's. posterior predictive distribution? First find it symbolically, then numerically.
- (f) If Eygor brought back 10 more brains from the local morgue, would Dr. F. be able to update her prior in some way that is more efficient than doing the individual update procedure 10 times?

7 Coding Question: Two Gaussians [25 Points]

We saw in lecture that if we know $p(X, Y)$ then we can derive an optimal decision rule, but we were sad to realize that we never really know $p(X, Y)$. One strategy for addressing this problem is to directly estimate $p(X, Y)$ and then use the estimate to derive a decision rule that would be optimal if our estimate were accurate.

In this question we consider a generative model for a dataset comprised of a mixture of two gaussians. The data is generated as follows. Let $C_0 = N(\mu_0, \Sigma_0)$ and $C_1 = N(\mu_1, \Sigma_1)$ be two gaussians where μ_0 and $\mu_1 \in \mathbb{R}^d$ are the means and Σ_0 and $\Sigma_1 \in \mathbb{R}^{d \times d}$ are two covariances. Let $y \in \{0, 1\}$ be a latent variable indicating if x is drawn from C_0 or C_1 . The probability density of x is defined as follows

$$P(x) = P(x|y = 1)P(y = 1) + P(x|y = 0)P(y = 0) \quad (1)$$

Our goal is to derive and implement the bayes optimal classifier δ such that given a new point $x' \in \mathbb{R}^d$,

$$\delta(x') = \arg \max_{y \in \{0,1\}} P(x'|y) \quad (2)$$

We have provided two csv files train.csv and test.csv for the completion of this question.

- (a) (Empirics) From train.csv, what is your maximum likelihood estimate for $P(y = 0)$ and $P(y = 1)$? What is your estimate for μ_0 and μ_1 ? What is your estimate for Σ_0 and Σ_1 ? Do you notice something about Σ_0 and Σ_1 ? (Hint: Don't overthink)
- (b) (Theory) What are $P(y = 1|x)$ and $P(y = 0|x)$ proportional to, as a function of x ?
- (c) (Theory) Derive an equation for the decision boundary for $x \in \mathbb{R}^d$ where

$$\ln(P(y = 1|x)) = \ln(P(y = 0|x)) \quad (3)$$

Here we compare the log likelihood as it simplifies the derivation. Is this decision boundary (as a function of x) linear, quadratic, etc.? How does the decision boundary simplify when $\Sigma_0 = \Sigma_1$?

- (d) (Empirics) Using the decision boundary derived in part (c), classify the points in test.csv as $y = 0$ or 1 . It suffices to write down the form of the decision boundary and associated decision rule.

8 Optional

8.1 Throwing rocks: asymmetric loss

You just bought a new trebuchet and you are interested in making predictions about how far it can throw a rock. Your ballistics officer tells you that optimizing the squared error of your predictions is not appropriate for the problem. If your prediction is within some constant c of the true value then they can use your predicted value to aim the trebuchet such that it hits the castle, but if your prediction off by more than c , then using the prediction for aiming will cause the trebuchet to miss. However, the fact is, it's better for your prediction to be too short than too far.

So, we will let

$$L(a, g) = \begin{cases} 0 & \text{if } |a - g| < c \\ 1 & \text{if } a - g > c \\ 2 & \text{if } g - a > c \end{cases}$$

If you know that the range of the ball for these types of trebuchets is distributed as a Gaussian with mean μ and variance σ^2 , what prediction minimizes loss L ?

(This is a little bit tricky. It's fine to just write out an expression in terms of c , μ , and σ^2 .)

8.2 Copy that: discrete Bayes update and decision theory

You have just bought a copy machine at a garage sale. You know it is one of two possible models, m_1 or m_2 , but the tag has fallen off, so you're not sure which.

You do know that m_1 machines have a 0.1 "error" (bad copy) rate and m_2 machines have a 0.2 error rate.

- (a) You use your machine to make 1000 copies, and 140 of them are bad. What is the maximum likelihood estimate of the machine's error rate? Explain why. (Remember that you're sure it's one of those two types of machines).
- (b) Looking more closely, you can see part of the label, and so you think that, just based on the label it has a probability 0.2 of being an m_1 type machine and a probability 0.8 of being an m_2 type machine. If you take that to be your prior, and incorporate the data from part a, what is your posterior distribution on the type of the machine?
- (c) Given that posterior, what is the probability that the next copy will be a failure?
- (d) You intend to sell this machine on the web. Because it's used, you have to sell it with a warranty. You can offer a gold or a silver warranty. If it has a gold warranty and the buyer runs it for 1000 copies and gets more than 150 bad copies, then you are obliged to pay \$1000 in damages; if it has a silver warranty, you have to pay damages if it generates more than 300 bad copies in 1000 copies. Your maximum reasonable asking price for a machine with a gold warranty is \$300; for a machine with a silver warranty, it is \$100. You can assume the machine will sell at these prices. What type of warranty should you offer on this machine?
- (e) Under what conditions would it be better to just throw the machine away, rather than try to sell it?

8.3 Dirichlet Priors

Exercise borrowed from Stat180 at UCLA. See Bishop, sections 2.1 and 2.2 for background on Beta and Dirichlet distributions.

The Dirichlet distribution is a multivariate version of the Beta distribution. When we have a coin with two outcomes, we really only need a single parameter θ to model the probability of heads. But now let's consider a "thick" coin that has three possible outcomes: heads, tails, and edge. Now we need two parameters: θ_h is the probability of heads, θ_t is the probability of tails, and then the probability of an edge is $1 - \theta_h - \theta_t$.

The random variables $(V, W) \in [0, 1]$ and such that $V + W \leq 1$ have a Dirichlet distribution with parameters $\alpha_1, \alpha_2, \alpha_3$ if their joint density is

$$f(v, w) = v^{\alpha_1-1} w^{\alpha_2-1} (1 - v - w)^{\alpha_3-1} \frac{\Gamma(\alpha_1 + \alpha_2 + \alpha_3)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\Gamma(\alpha_3)} .$$

This is a direct generalization of the Beta distribution. (Note that Γ refers to the Gamma function, which is a generalization of factorial.)

- (a) If (θ_h, θ_t) have a Dirichlet distribution as above, what is the marginal distribution of θ_h ?
- (b) Suppose you are playing with a thick coin, and get results $x^{(1)} \dots x^{(n)}$, resulting in H heads and T tails out of n throws. Given θ_h and θ_t the random variables H and T have a multinomial distribution:

$$\Pr(H, T | \theta_h, \theta_t) = \frac{n!}{H!T!(n-H-T)!} \theta_h^H \theta_t^T (1 - \theta_h - \theta_t)^{n-H-T} .$$

Assume a uniform prior on the space of possible values of θ_h and θ_t (remembering that they are constrained such that $\theta_h \geq 0$, $\theta_t \geq 0$, and $\theta_h + \theta_t \leq 1$). What is the posterior distribution for θ_h and θ_t ?

- (c) In this same setting, what is the predictive distribution for getting another head? That is, what's $\Pr(x^{(n+1)} = \text{heads} \mid x^{(1)} \dots x^{(n)})$?
- (d) Now assume a Dirichlet prior for θ_h and θ_t with parameters $\alpha_1, \alpha_2, \alpha_3$. What is the posterior in this case?
- (e) In this same case, what is the predictive distribution?
- (f) If you assume a squared-error loss on the predicted parameter, that is,

$$L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2 ,$$

what is the Bayes-optimal estimate of θ_h and θ_t ?

- (g) As $n \rightarrow \infty$, how do optimal estimates relate to the maximum likelihood estimates and to the prior?

8.4 More fun with Dirichlet

Given a parameterized family of probability models $\Pr(x \mid \theta)$ and a data set $D = (x^{(1)}, \dots, x^{(n)})$ comprised of independent samples $x^{(i)} \approx \Pr(x \mid \theta)$, we fit the model to the data so as to maximize the likelihood (or log-likelihood) of all samples. This gives the maximum-likelihood (ML) estimate of the parameters:

$$\hat{\theta}_{ML} = \arg \max_{\theta} \log \Pr(D \mid \theta)$$

This approach does not express any prior bias as to which values of θ we should prefer when data is limited.

In the sequel, we consider a regularized approach to parameter estimation. Here, we specify a prior model $\Pr(\theta)$ over the set of allowed parameter settings Θ . Given a prior model, we may then employ Bayes' rule to compute the posterior probability of θ given the observations:

$$\Pr(\theta \mid D) = \frac{\Pr(D \mid \theta) \Pr(\theta)}{\Pr(D)}$$

where

$$\Pr(D) = \int_{\Theta} \Pr(D \mid \theta) \Pr(\theta) d\theta$$

Then, we fit the model to the data by maximizing the (log-) probability of θ conditioned on the data,

$$\begin{aligned} \hat{\theta}_{MAP} &= \arg \max_{\theta} \log \Pr(\theta \mid D) \\ &= \arg \max_{\theta} \{\log \Pr(D \mid \theta) + \log \Pr(\theta) - \log \Pr(D)\} \\ &= \arg \max_{\theta} \{\log \Pr(D \mid \theta) + \log \Pr(\theta)\} \end{aligned}$$

Note that we have dropped the $-\log \Pr(D)$ term as this does not depend upon θ and does not affect the parameter estimate. Hence, we do not need to explicitly evaluate the integral in the denominator. This may be viewed as a penalized log-likelihood criterion, i.e. we maximize $J(\theta) = \log \Pr(D; \theta) + f(\theta)$ subject to the regularization penalty $f(\theta) = \log \Pr(\theta)$. The parameter estimate $\hat{\theta}_{\text{MAP}}$ is known as the maximum a posteriori (MAP) estimate.

In this problem you will construct MAP estimates for the probabilities of a (potentially biased) M -sided die, i.e. $x^{(i)} \in \{1, \dots, M\}$. We consider the fully-parameterized representation $\Pr(x = k) = \theta_k$, where $0 \leq \theta_k \leq 1$ for $k = 1, \dots, M$ and $\sum_k \theta_k = 1$. This simple model has many relevant applications.

Consider a document classification task, where we need class-conditional distributions over words in the documents. Suppose we only consider words $1, \dots, M$ (for relatively large M). Each word in the document is assumed to have been drawn at random from the distribution $\Pr(x = k | y; \theta) = \theta_{k|y}$, where $\sum_k \theta_{k|y} = 1$ for each class y . Thus the selection of words according to the distribution $\theta_{k|y}$ can be interpreted as a (biased) M -sided die.

Now, the probability of generating all words $x^{(1)}, \dots, x^{(n)}$ in a document of length n would be

$$\Pr(D | y; \theta) = \prod_{i=1}^n \Pr(x^{(i)} | y; \theta) = \prod_{i=1}^n \theta_{x^{(i)}|y}$$

assuming the document belongs to class y . Note that this model cares about how many times each word occurs in the document. It is a valid probability model over the set of words in the document.

Since we typically have very few documents per class, it is important to regularize the parameters, i.e., provide a meaningful prior answer to the class conditional distributions.

Let's start by briefly revisiting ML estimation of the (biased) M -sided die. Similarly to calculations you have already performed, the ML estimate of the parameter θ from n samples is given by the empirical distribution:

$$\hat{\theta}_x = \frac{n(x)}{n}$$

where $n(x)$ is the number of times value x occurred in n samples. The count $n(x)$ is also a *sufficient statistic* for θ_x as it is all we need to know from the available n samples in order to estimate θ_x .

Next, we consider MAP estimation. To do so, we must introduce a prior distribution over the θ 's. A natural choice for this problem is the Dirichlet distribution

$$\Pr(\theta; \beta) = \frac{1}{Z(\beta)} \prod_{k=1}^M \theta_k^{\beta_k}$$

with non-negative hyperparameters $\beta = (\beta_k > 0, k = 1, \dots, M)$ and where $Z(\beta)$ is just the normalization constant (which you saw earlier and which you do not need to evaluate in this problem).

- (a) First, consider this prior model (ignoring the data for the moment). What value of θ is most likely under this prior model? That is, compute

$$\hat{\theta}(\beta) = \arg \max_{\theta} \log \Pr(\theta; \beta)$$

This is the *a priori* estimate of θ before observing any data.

- (b) Next, given the data D , compute the MAP estimate of θ as a function of the hyperparameters β and the data D (use the sufficient statistics $n(x)$):

$$\hat{\theta}_{\text{MAP}}(D; \beta) = \arg \max_{\theta} \log \Pr(\theta | D; \beta)$$

Note that you do not need to calculate $Z(\beta)$ in order to perform this optimization; you can optimize the penalized log-likelihood $J(\theta) = \log \Pr(D | \theta) + f(\theta; \beta)$ with a simple penalty function $f(\theta; \beta)$, as discussed above. Thus we do not have to evaluate the full posterior distribution $\Pr(\theta | D; \beta)$ in order to perform the regularization.

- (c) Show that your MAP estimate may be expressed as a convex combination of the a priori estimate $\hat{\theta}(\beta)$ and the ML estimate $\hat{\theta}_{\text{ML}}(D)$. The means that we may write

$$\hat{\theta}_{\text{MAP}}(D; \beta) = (1 - \lambda)\hat{\theta}_{\text{ML}}(D) + \lambda\hat{\theta}(\beta)$$

for some $\lambda \in [0, 1]$. Note that the same convex combination holds for each component θ_x . Determine λ as a function of the number of samples n and the hyperparameters β .

As this shows, one way of thinking of a prior distribution is that it is a proxy for any data we have observed in the past but no longer have available. The normalized parameters $\hat{\beta}_i = \beta_i/N$, where $N = \sum_i \beta_i$, express our prior estimate of the parameters θ while the normalization parameter N expresses how strongly we believe in that prior estimate.