

Bayesian Estimation & Conditional Models

CMPUT 296: Basics of Machine Learning

Textbook §5.2-5.3

Logistics

- **Assignment 1 is due TODAY at 11:59pm**
- Assignment 2 will be available tomorrow
- Thought Questions #2: Ch.4-6 are due **October 8**

Recap: Parameter Estimation

- We are usually interested in predicting the value of unseen data X_{n+1} based on **training data** $\mathcal{D} = \{x_1, \dots, x_n\}$
- Instead, we will want to choose a **model** \hat{f} from a **hypothesis space** \mathcal{F}
 - Where the data are generated according to some "true" model f^*
 - \mathcal{F} is often **parametric**: its members identified by **parameter** values
 - So choosing a **model** is equivalent to choosing a set of **parameter** values
- Two approaches to parameter estimation:

$$f_{\text{MAP}} = \arg \max_{f \in \mathcal{F}} p(f \mid \mathcal{D}) = \arg \max_{f \in \mathcal{F}} p(\mathcal{D} \mid f)p(f)$$

$$f_{\text{MLE}} = \arg \max_{f \in \mathcal{F}} p(\mathcal{D} \mid f)$$

Outline

1. Recap & Logistics
2. Bayesian Estimation
3. Conditional Distributions

Point Estimates

- Suppose we have a dataset \mathcal{D} that was generated by a model $f(\cdot | \theta^*) \in \mathcal{F} = \{f(\cdot | \theta) | \theta \in \mathbb{R}\}$
- A **point estimate** asks: What is the **best single guess** for the parameter?
 - **MLE:** $\arg \max_{\theta} p(\mathcal{D} | \theta)$
 - **MAP:** $\arg \max_{\theta} p(\theta | \mathcal{D})$
 - Estimate of θ that has lowest **expected error**?

Bayes Estimator

The **Bayes estimator** is the point estimate that **minimizes the posterior risk** $c(\hat{\theta})$, where

$$c(\hat{\theta}) = \int_{\mathcal{F}} \ell(\theta, \hat{\theta}) p(\theta | \mathcal{D}) d\theta$$

The **loss** $\ell(\theta, \hat{\theta})$ expresses how "wrong" we are if we estimate $\hat{\theta}$ when the true answer is θ .

Bayes Estimator for Squared Loss

When $\ell(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$:

$$c(\hat{\theta}) = \int_{\mathcal{F}} (\theta - \hat{\theta})^2 p(\theta | \mathcal{D}) d\theta$$

$$\iff \frac{\partial}{\partial \hat{\theta}} c(\hat{\theta}) = \frac{\partial}{\partial \hat{\theta}} \int_{\mathcal{F}} (\theta - \hat{\theta})^2 p(\theta | \mathcal{D}) d\theta$$

$$= \int_{\mathcal{F}} \frac{\partial}{\partial \hat{\theta}} (\theta - \hat{\theta})^2 p(\theta | \mathcal{D}) d\theta$$

$$= 2 \int_{\mathcal{F}} (\hat{\theta} - \theta) p(\theta | \mathcal{D}) d\theta$$

$$= 2\hat{\theta} \int_{\mathcal{F}} p(\theta | \mathcal{D}) d\theta - 2 \int_{\mathcal{F}} \theta p(\theta | \mathcal{D}) d\theta$$

$$= 2\hat{\theta} - 2 \int_{\mathcal{F}} \theta p(\theta | \mathcal{D}) d\theta$$

$$\frac{\partial}{\partial \hat{\theta}} \ell(\hat{\theta}) = 0$$

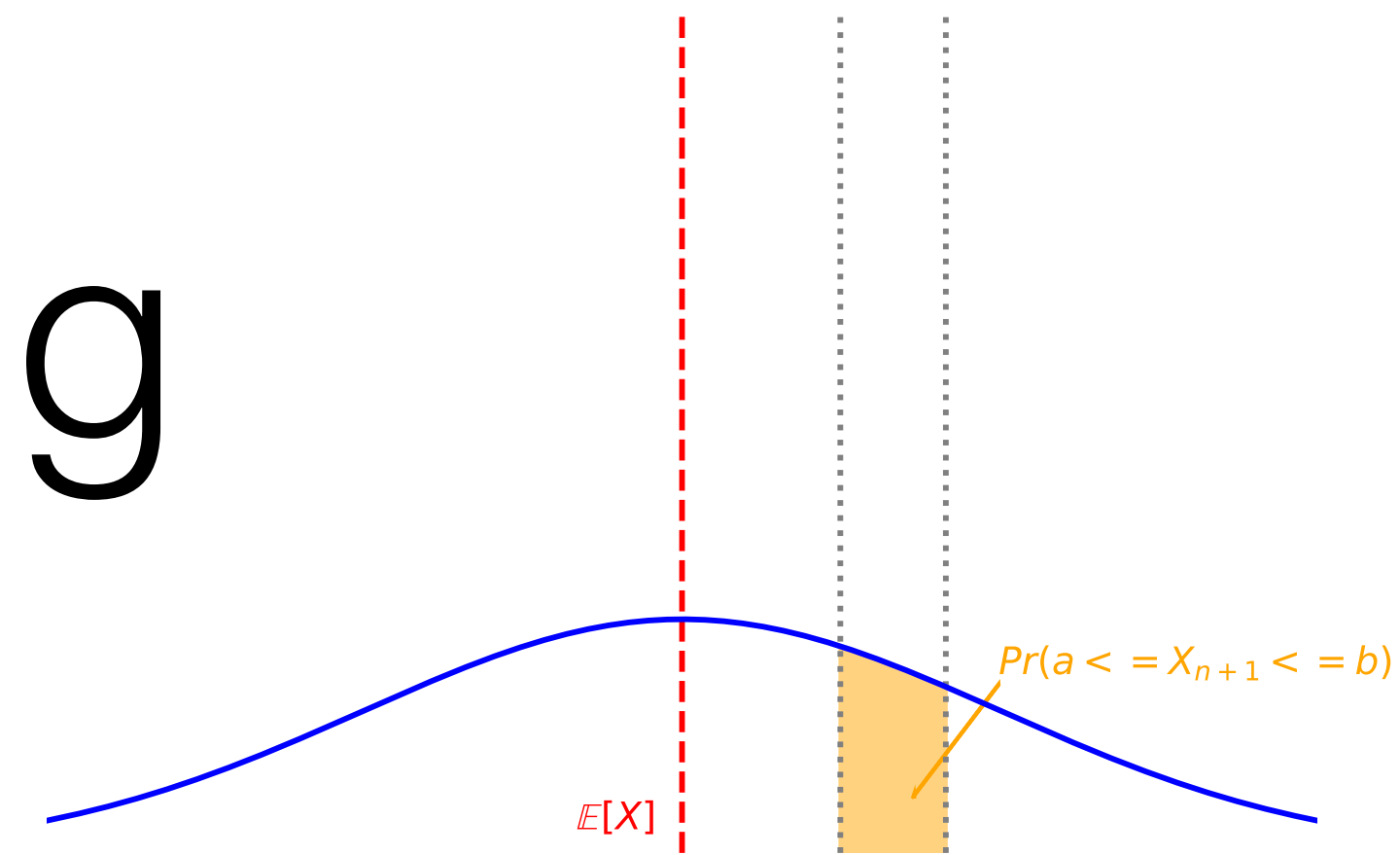
$$\iff 0 = 2\hat{\theta} - 2 \int_{\mathcal{F}} \theta p(\theta | \mathcal{D}) d\theta$$

$$\iff \hat{\theta} = \int_{\mathcal{F}} \theta p(\theta | \mathcal{D}) d\theta$$

$$= \mathbb{E}[\theta | \mathcal{D}]$$

Bayesian Reasoning

Question: How can we answer $\Pr(a \leq X_{n+1} \leq b)$?



1. **MLE:** $F(b | \theta_{\text{MLE}}) - F(a | \theta_{\text{MLE}})$

2. **MAP:** $F(b | \theta_{\text{MAP}}) - F(a | \theta_{\text{MAP}})$

3. **Bayes optimal estimator:** $F(b | \theta_B) - F(a | \theta_B)$ ←

Question: Does this use of θ_B make sense? Why?

4. **Bayesian:** $\int_{\mathcal{F}} [F(b | \theta) - F(a | \theta)] p(\theta | \mathcal{D}) d\theta$

$$= \mathbb{E} [F(b | \theta) - F(a | \theta) | \mathcal{D}]$$

Example: Poisson Data with Gamma Prior

Example: Suppose dataset $\mathcal{D} = \{8, 4, 5, 9, 5, 2\}$ is drawn i.i.d. from an unknown Poisson distribution, with parameter λ_0 . We have a Gamma prior over λ ; that is,

$$\text{prior } p(\lambda) = \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)}$$

$$\text{and likelihood } p(\mathcal{D} | \lambda) = \frac{\lambda^{\left(\sum_{i=1}^n x_i\right)} e^{-n\lambda}}{\prod_{i=1}^n x_i!}.$$

To compute the Bayes estimator, we will need the full **posterior** $p(\lambda | \mathcal{D})$, and not just the joint $p(\mathcal{D} | \lambda)p(\lambda) \propto p(\lambda | \mathcal{D})$. (**Why?**)

That means we need to compute the **model evidence** as well.

Poisson Data with Gamma Prior 2

$$\begin{aligned}
 p(\mathcal{D}) &= \int_0^{\infty} p(\mathcal{D} | \lambda) p(\lambda) d\lambda \\
 &= \int_0^{\infty} \frac{\lambda^{\left(\sum_{i=1}^n x_i\right)} e^{-n\lambda}}{\prod_{i=1}^n x_i!} \cdot \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)} d\lambda \\
 &= \frac{\Gamma\left(k + \sum_{i=1}^n x_i\right)}{\theta^k \Gamma(k) \prod_{i=1}^n x_i! \left(n + \frac{1}{\theta}\right)^{\left(k + \sum_{i=1}^n x_i\right)}}
 \end{aligned}$$

$$\begin{aligned}
 p(\lambda | \mathcal{D}) &= \frac{\lambda^{\left(\sum_{i=1}^n x_i\right)} e^{-n\lambda}}{\prod_{i=1}^n x_i!} \cdot \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)} \cdot \frac{\theta^k \Gamma(k) \prod_{i=1}^n x_i! \left(n + \frac{1}{\theta}\right)^{\left(k + \sum_{i=1}^n x_i\right)}}{\Gamma\left(k + \sum_{i=1}^n x_i\right)} \\
 &= \frac{\lambda^{\left(k + \sum_{i=1}^n x_i\right) - 1} \cdot e^{-\lambda\left(n + \frac{1}{\theta}\right)} \cdot \left(n + \frac{1}{\theta}\right)^{\left(k + \sum_{i=1}^n x_i\right)}}{\Gamma\left(k + \sum_{i=1}^n x_i\right)} \\
 &= \frac{\lambda^{\left(k + \sum_{i=1}^n x_i\right) - 1} \cdot e^{-\lambda\left(n + \frac{1}{\theta}\right)}}{\left(\frac{1}{n + \frac{1}{\theta}}\right)^{\left(k + \sum_{i=1}^n x_i\right)} \cdot \Gamma\left(k + \sum_{i=1}^n x_i\right)}
 \end{aligned}$$

i.e., a **Gamma distribution** with

$$k' = k + \sum_{i=1}^n x_i \text{ and } \theta' = \frac{\theta}{n\theta + 1} = \frac{1}{n + 1/\theta}$$

Fun fact: $\int_0^{\infty} x^{\alpha-1} e^{-\beta x} dx = \frac{\Gamma(\alpha)}{\beta^\alpha}$

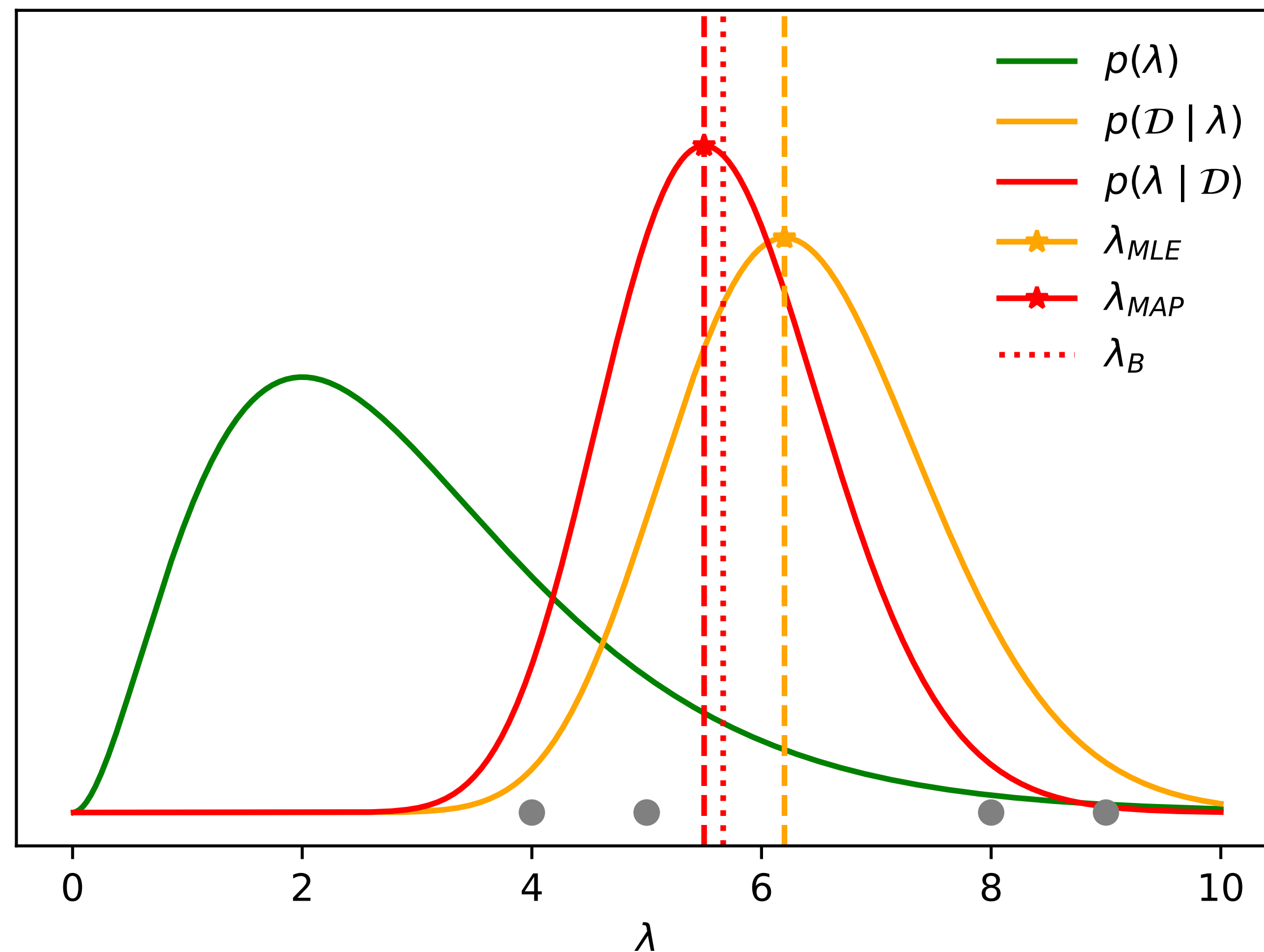
Conjugate Priors

- **Gamma** is a **conjugate prior** for the **Poisson** distribution
- Starting from **prior** $\text{Gamma}(k, \theta)$ and assuming a Poisson **likelihood**, after seeing x_1, \dots, x_n , **posterior** is $\text{Gamma}\left(k + \sum_{i=1}^n s_i, \frac{1}{n + 1/\theta}\right)$
- Similarly, **Beta** is a conjugate prior for the **Binomial** distribution
- Starting from **prior** $\text{Beta}(a, b)$ and assuming a Binomial **likelihood**, after seeing n_T successes and n_F failures, **posterior** is $\text{Beta}(a + n_T, b + n_F)$.

Poisson Data Example: Updating

Example: Suppose dataset $\mathcal{D} = \{8, 4, 5, 9, 5, 2\}$ is drawn i.i.d. from an unknown Poisson distribution, with parameter λ_0 with prior **Gamma**($k = 3, \theta = 1$):

$$\mathcal{D} = \{8, 4, 5, 9, 5\}$$



Advanced: Bayesian Methods with Nonconjugate Priors

- Conjugate priors are very convenient, and you should use them wherever possible
- **Question:** What can we do if the priors are **not conjugate**?

- In general, the integral to compute $p(\mathcal{D})$ will be **intractable**

- The usual technique is variants of **Monte Carlo sampling**

- *Basic idea:* Generate a bunch of random samples $\theta_1, \dots, \theta_R \stackrel{i.i.d}{\sim} p(\theta | \mathcal{D})$

$$\mathbb{E} \left[\frac{1}{R} \sum_{r=1}^R f(\theta_r) \right] = \mathbb{E}[f(\theta)]$$

- There are multiple techniques for generating random samples from an **unnormalized distribution** $q(\theta) \propto p(\theta)$

- We can use one of these techniques to sample from $p(\mathcal{D} | \theta)p(\theta) \propto p(\theta | \mathcal{D})$

Conditional Models

Question: How can we ask "With what probability is  an image of a cat" using the models we have been learning?

- We want a **different distribution** depending on the **image**, and
- We want to be able to ask about **multiple images**
- Given an image described by pixels \mathbf{x} , we want something like

$$\Pr(Y = \text{cat} \mid X = \mathbf{x})$$

- Our models can be parameterized families of **conditional distributions**:

$$\mathcal{F} = \{f(y, x; \theta) \mid \theta \in \mathbb{R}^k\}$$

MLE, MAP, Bayesian Prediction for Conditional Models

Given a hypothesis space $\mathcal{F} = \{p(\cdot | \cdot, \theta) \mid \theta \in \mathbb{R}\}$ and a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ of observed **features** x_i and their corresponding **labels** y_i :

MLE: $p(y | x) = p(y | x, \theta_{\text{MLE}})$ where $\theta_{\text{MLE}} = \arg \max_{\theta} \sum_i \ln p(y_i | x_i, \theta)$

MAP: $p(y | x) = p(y | x, \theta_{\text{MLE}})$ where $\theta_{\text{MAP}} = \arg \max_{\theta} \ln p(\theta) + \sum_i \ln p(y_i | x_i, \theta)$

Bayesian: $p(y | x) = \int_{\mathcal{F}} p(y | x, \theta) p(\theta | \mathcal{D}) d\theta$

Question: What happened to θ_B ?

Summary

- The MAP, MLE, and Bayes estimators for a model parameter are all **point estimates**
- MAP and MLE can be computed without computing $p(\mathcal{D})$
- **Conjugate priors** make it possible to perform Bayesian updates analytically
 - But many models don't have conjugate priors
- **Conditional models** allow us to change predictions based on **input features**
 - **MAP, MLE:** simply plug the features into the point estimate model
 - **Bayesian:** take posterior expected prediction over all models