

CS5340

Uncertainty Modeling in AI

Lecture 2: Fitting Probability Models

Asst. Prof. Harold Soh

AY 22/23

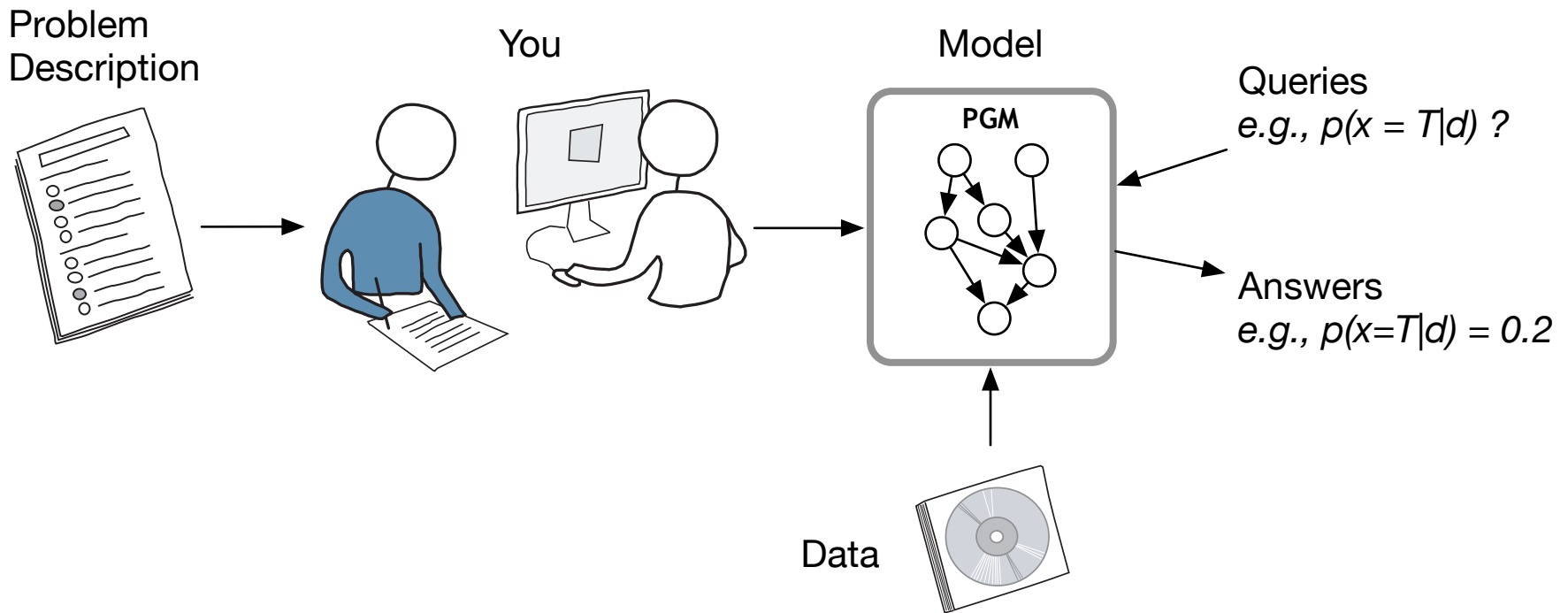
Semester 2

Course Schedule

Week	Date	Lecture Topic	Tutorial Topic
1	12 Jan	Introduction to Uncertainty Modeling + Probability Basics	Introduction
2	19 Jan	Simple Probabilistic Models	Probability Basics
3	26 Jan	Bayesian networks (Directed graphical models)	More Basic Probability
4	2 Feb	Markov random Fields (Undirected graphical models)	DGM modelling and d-separation
5	9 Feb	Variable elimination and belief propagation	MRF + Sum/Max Product
6	16 Feb	Factor graph and the junction tree algorithm	Quiz 1
-	-	RECESS WEEK	
7	2 Mar	Mixture Models and Expectation Maximization (EM)	Linear Gaussian Models
8	9 Mar	Hidden Markov Models (HMM)	Probabilistic PCA
9	16 Mar	Monte-Carlo Inference (Sampling)	Linear Gaussian Dynamical System
10	23 Mar	Variational Inference	MCMC + Sequential VAE
11	30 Mar	Inference and Decision-Making (Special Topic)	Quiz 2
12	6 Apr	Gaussian Processes (Special Topic)	Wellness Day
13	13 Apr	Project Presentations	Closing

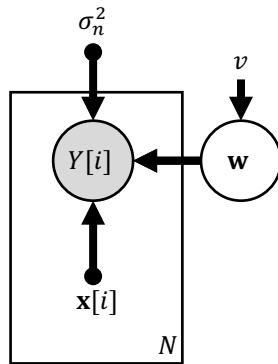
CS5340 in a nutshell

CS5340 is about how to “**represent**” and “**reason**”
with **uncertainty** in a computer.



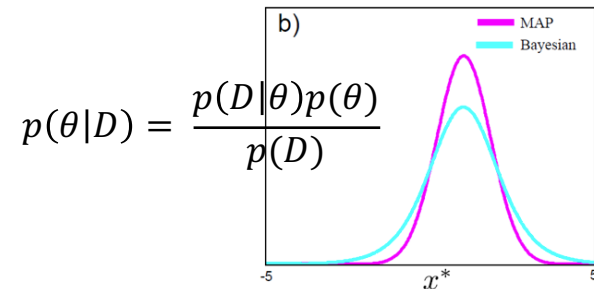
CS5340 in a nutshell

CS5340 is about how to “**represent**” and “**reason**” with **uncertainty** in a computer.



Representation: The *language* is probability and probabilistic graphical models (PGM).

The language is used to **model problems**.

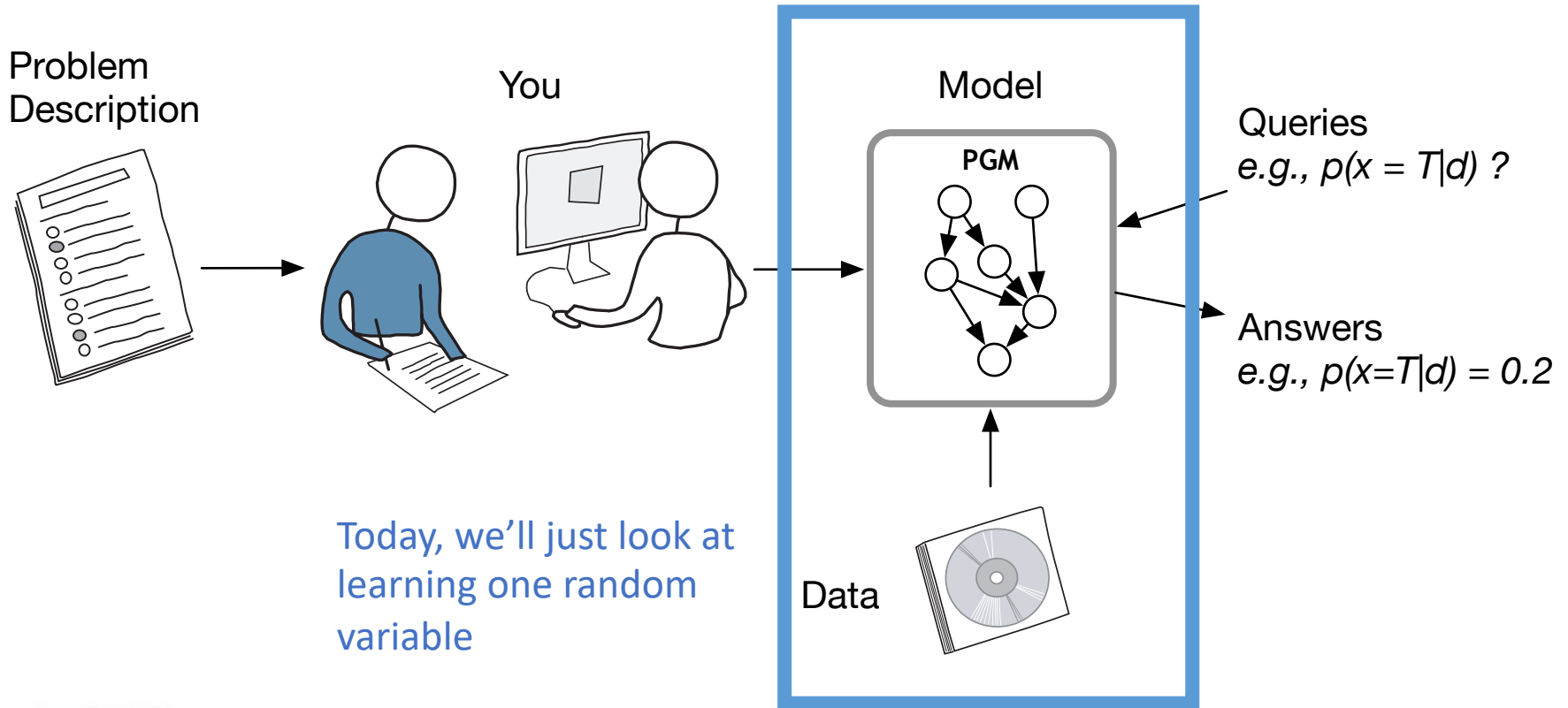


Reasoning: We use learning and inference algorithms to answer questions.

e.g., Belief-propagation/sum-product, MCMC, and variational Bayes

CS5340 in a nutshell

CS5340 is about how to “**represent**” and “**reason**” with **uncertainty** in a computer.



Summary: Sum and Product Rules

- Sum rule:

$$p(x) = \int p(x, y) dy$$

$$p(x) = \sum_y p(x, y)$$

- Product/Chain rule:

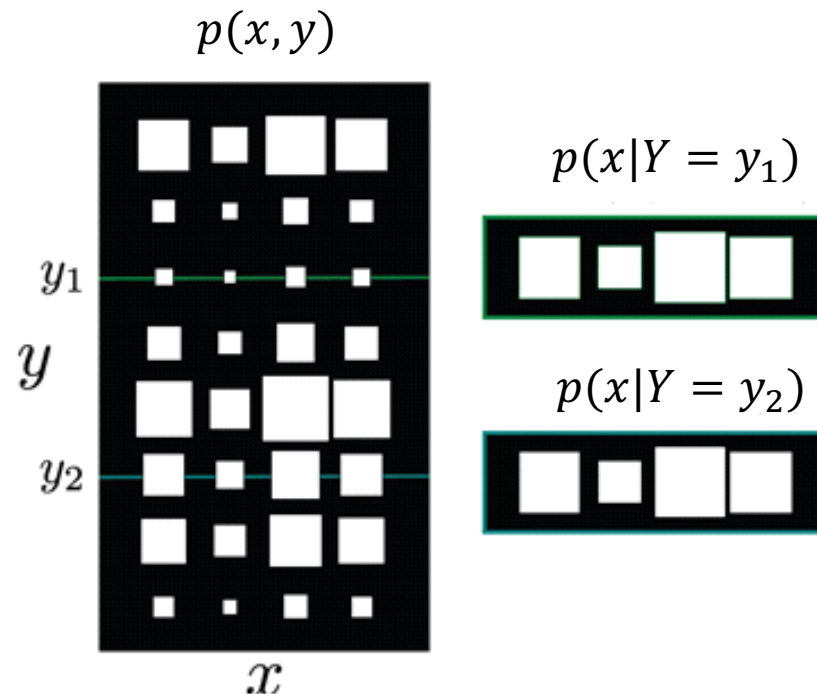
$$p(x, y) = p(x|y)p(y)$$

Probability: Independence

- The independence of X and Y means that **every conditional distribution is the same**.
- The value of Y **tells us nothing** about X and vice-versa.

$$p(x|y) = p(x)$$

$$p(y|x) = p(y)$$



Images Source: "Computer Vision: Models, Learning, and Inference", Simon Prince

Probability: Bayes' Rule

- Recall:

$$p(x, y) = p(x|y)p(y)$$
$$p(x, y) = p(y|x)p(x)$$



Thomas Bayes
1701–1761

- Eliminating $p(x, y)$, we get:

$$p(y|x)p(x) = p(x|y)p(y)$$

- Rearranging:

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\int p(x, y)dy} = \frac{p(x|y)p(y)}{\int p(x|y)p(y)dy}$$

Image source: "Pattern Recognition and Machine Learning", Christopher Bishop

Probability: Expectation

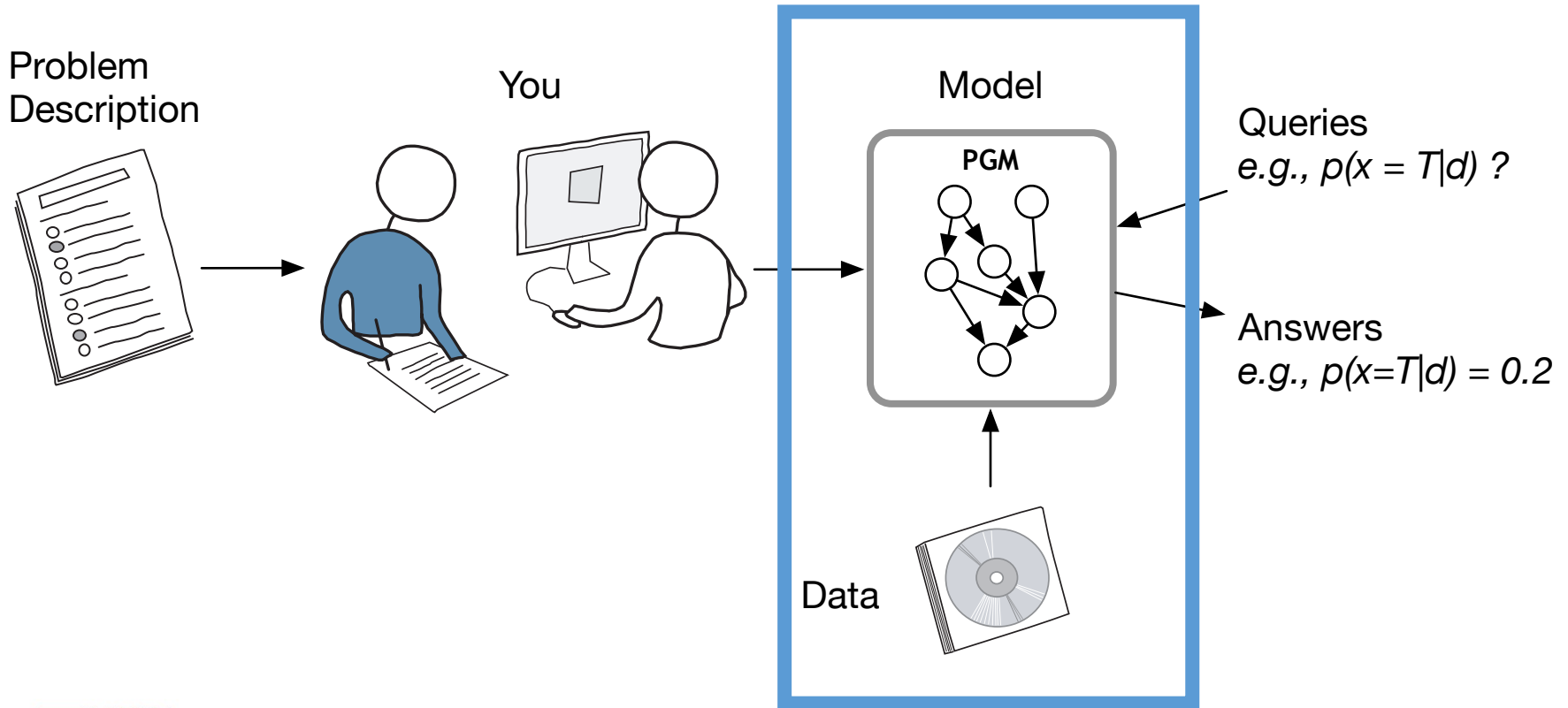
- The **expected or average value** of some function $f[x]$ taking into account the distribution of X .

Definition:

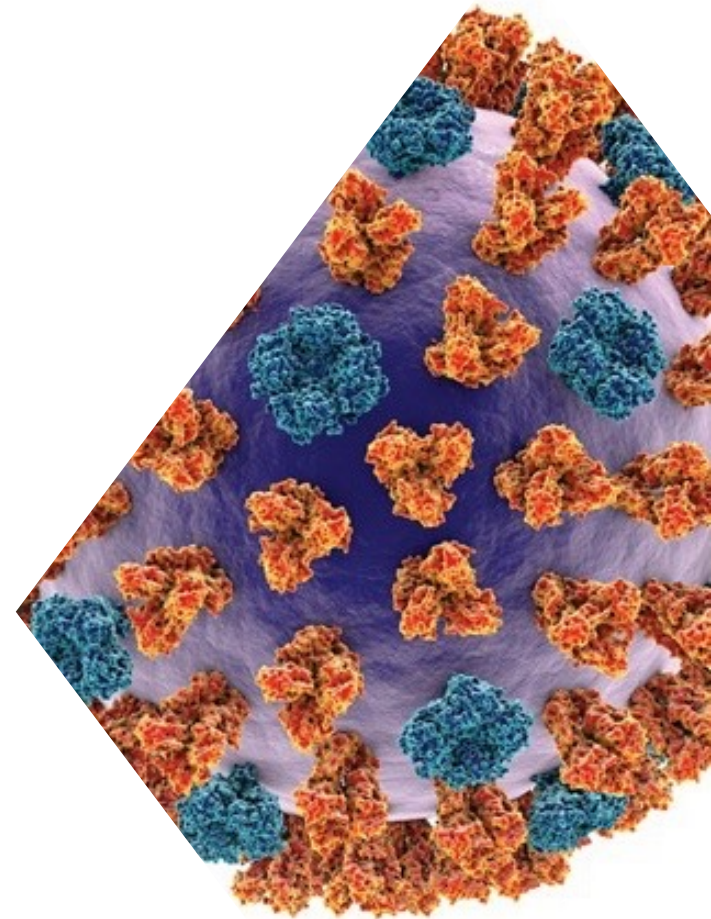
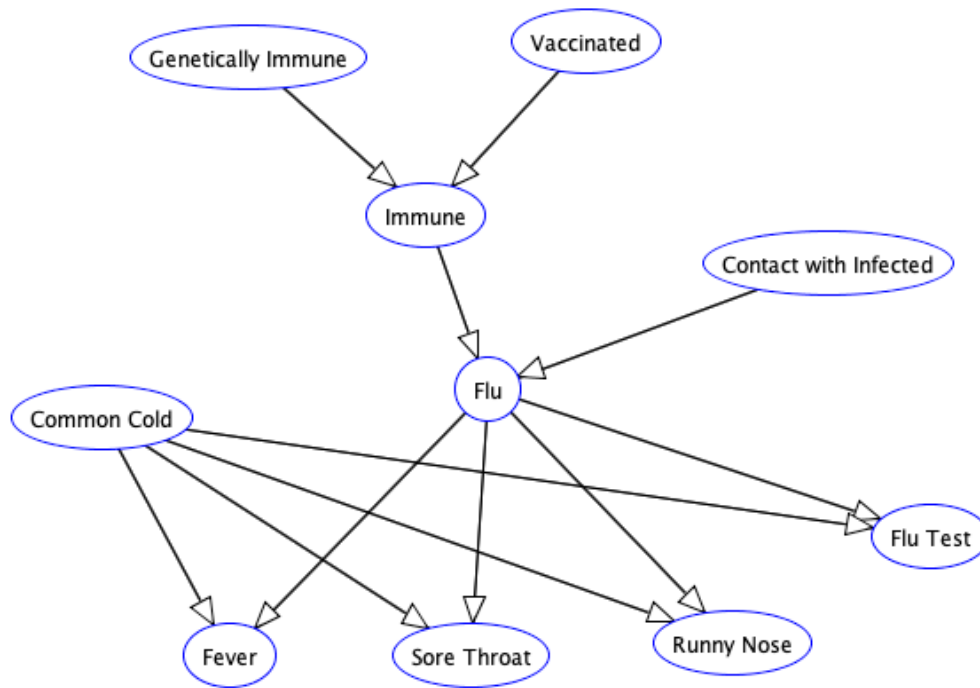
$$E[f[x]] = \sum_x f[x]p(x)$$
$$E[f[x]] = \int f[x]p(x)dx$$

CS5340 in a nutshell

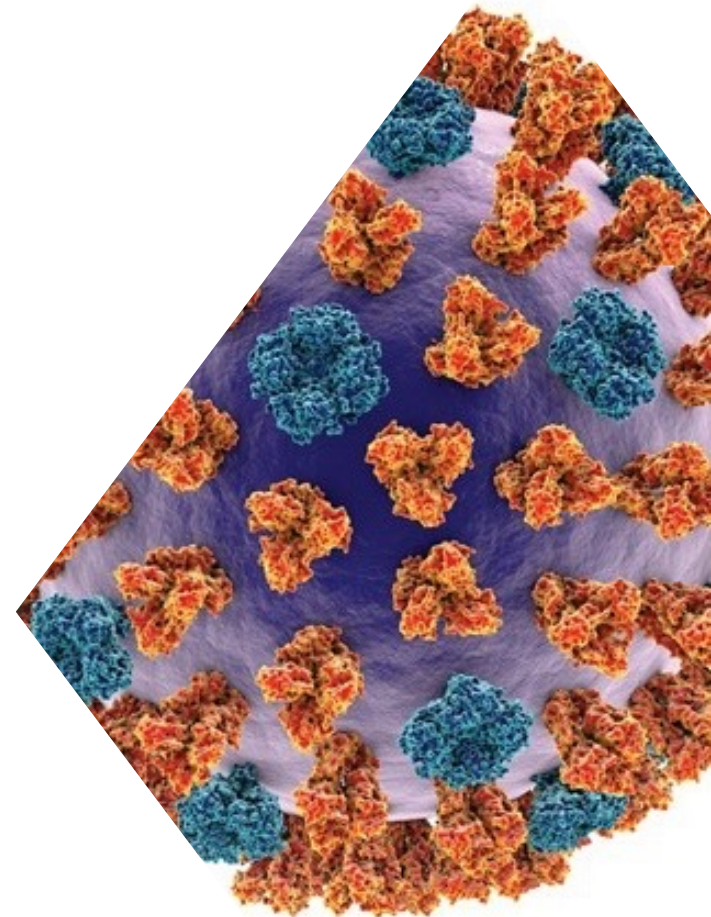
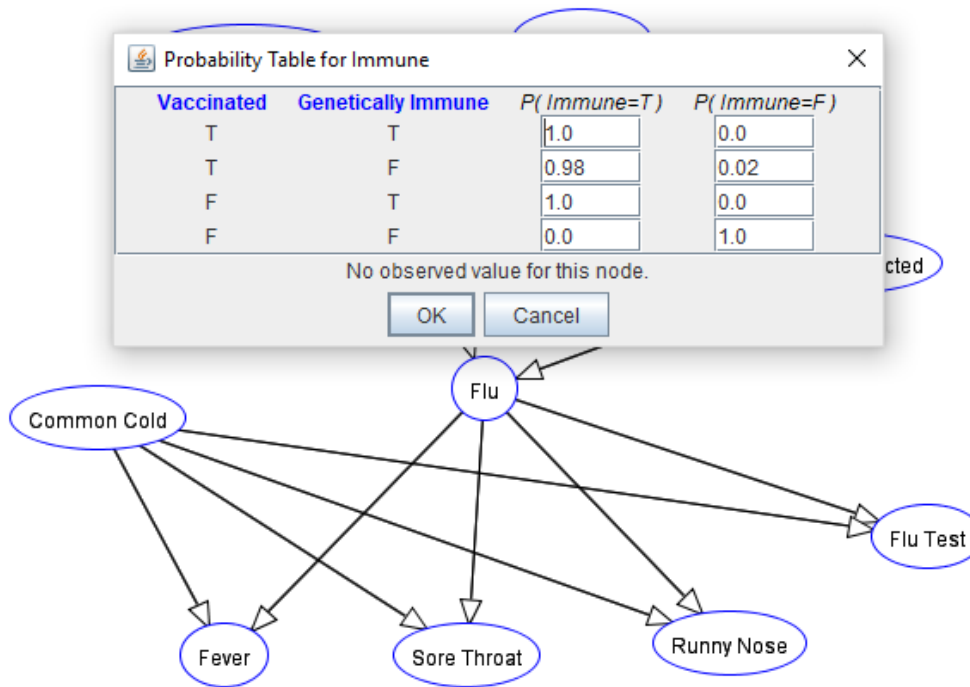
CS5340 is about how to “**represent**” and “**reason**”
with **uncertainty** in a computer.



Generative (Causal) Modeling of Relationships between Variables



Generative (Causal) Modeling of Relationships between Variables



Fitting Probability Models

- Focus on **parametric probability distributions** $p(x|\theta)$.
- How to **learn the unknown parameters θ** from a set of given data, i.e. instances of the random variable, $\mathcal{D} = \{x[1], \dots, x[N]\}$.
- And then **use those parameters to make predictions**.

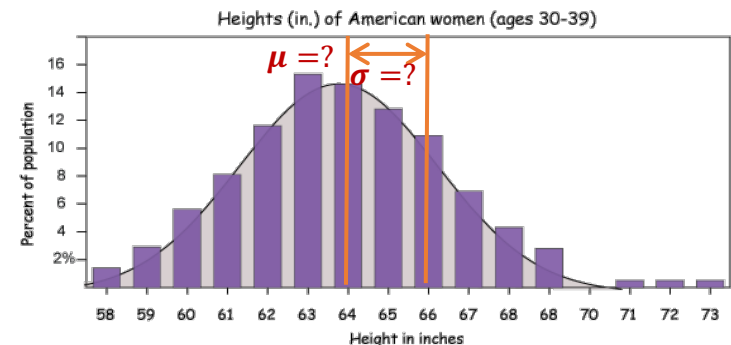


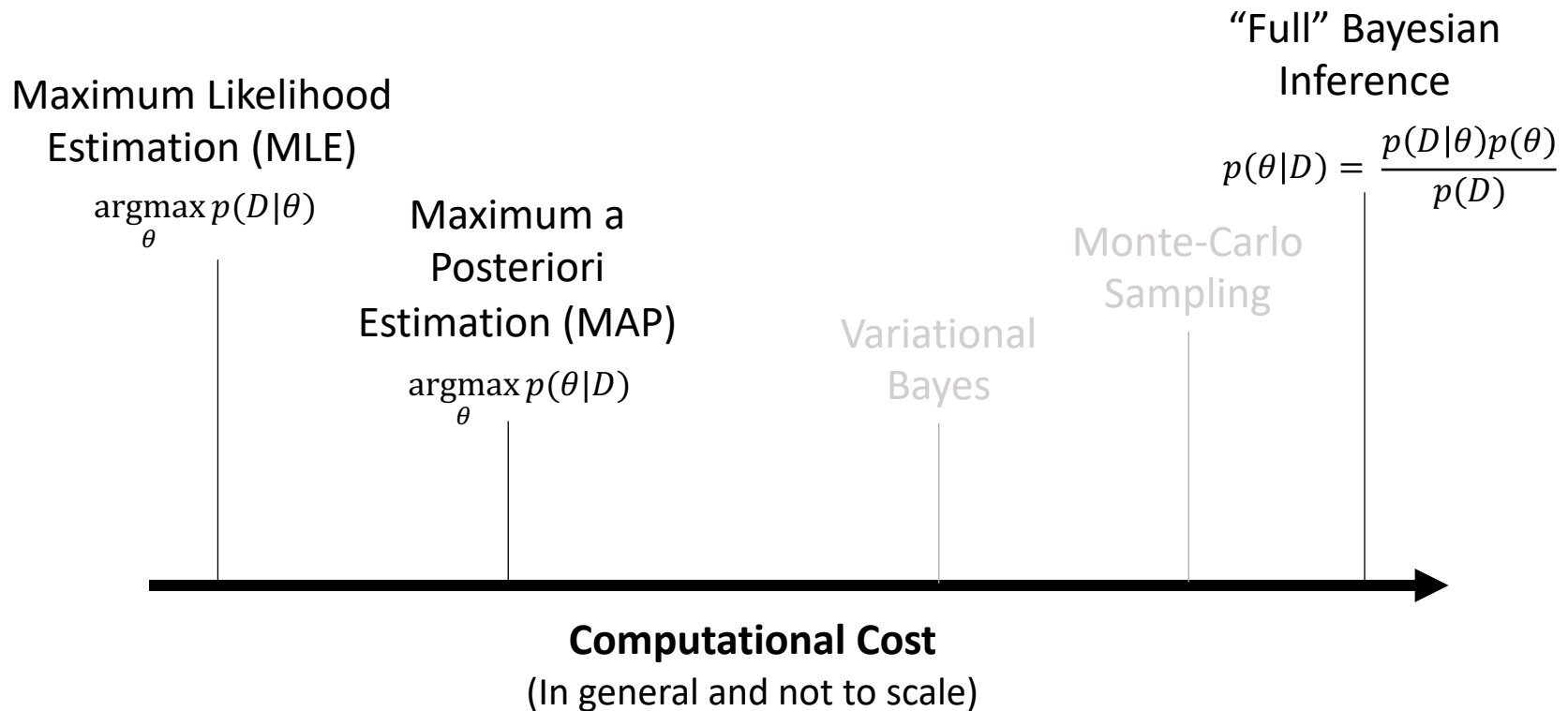
Image source: http://www.drcruzan.com/ProbStat_Distributions.html

Learning Outcomes

- Students should be able to:
 1. Use the **Maximum Likelihood**, **Maximum a Posteriori** and **Bayesian** approaches to learn the unknown parameters of probability distributions of a **single random variable** from data.
 2. Apply the assumption **independent and identically distributed samples** to simplify the parameter learning process.
 3. Apply the learned parameters to **make predictions**.
 4. Describe the **exponential family** and its properties

Learning Parameters

- Common approaches to **learn the unknown parameters** θ from a set of given data $\mathcal{D} = \{x[1], \dots, x[N]\}$:



Acknowledgements

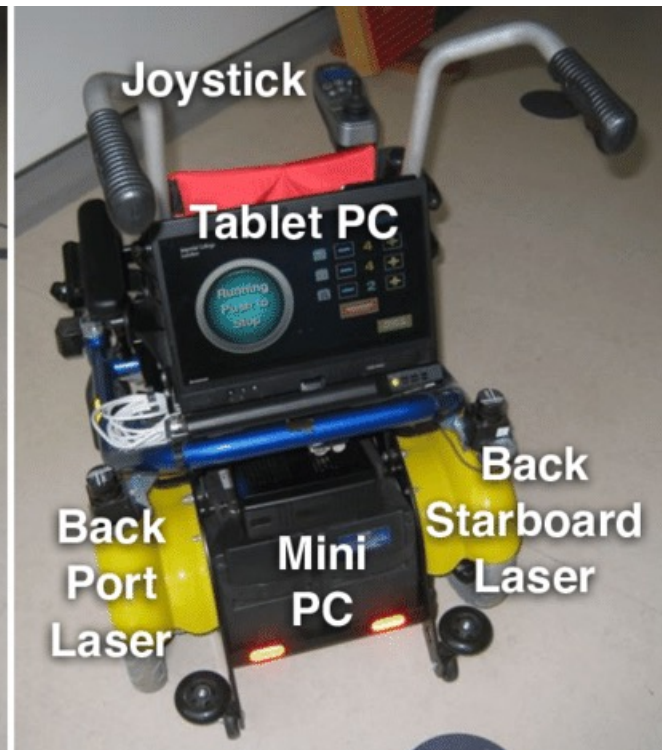
- A lot of slides and content of this lecture are adopted from:
 1. “Pattern Recognition and Machine Learning”, Christopher Bishop.
 2. “Computer Vision: Models, Learning, and Inference”, Simon Prince.
 3. Lee Gim Hee’s CS5340 slides.

Learning via MLE

Maximum Likelihood Estimation (MLE)

Building a *Smart* Wheelchair for Kids with Disabilities

- <https://youtu.be/XbyqU88jmb0>





smart mobility **ARTY** for kids

Problem: Sensor Uncertainty

- You have a ultrasonic ranger for your robot.
- Like other sensors, there is some error.
- How can you **model** and **estimate** the **uncertainty** of your range readings?
- *Later:* can we **predict the range** given a noisy reading?

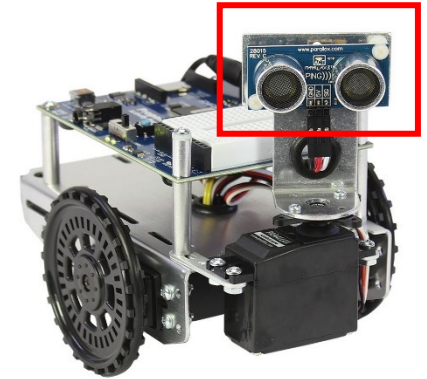
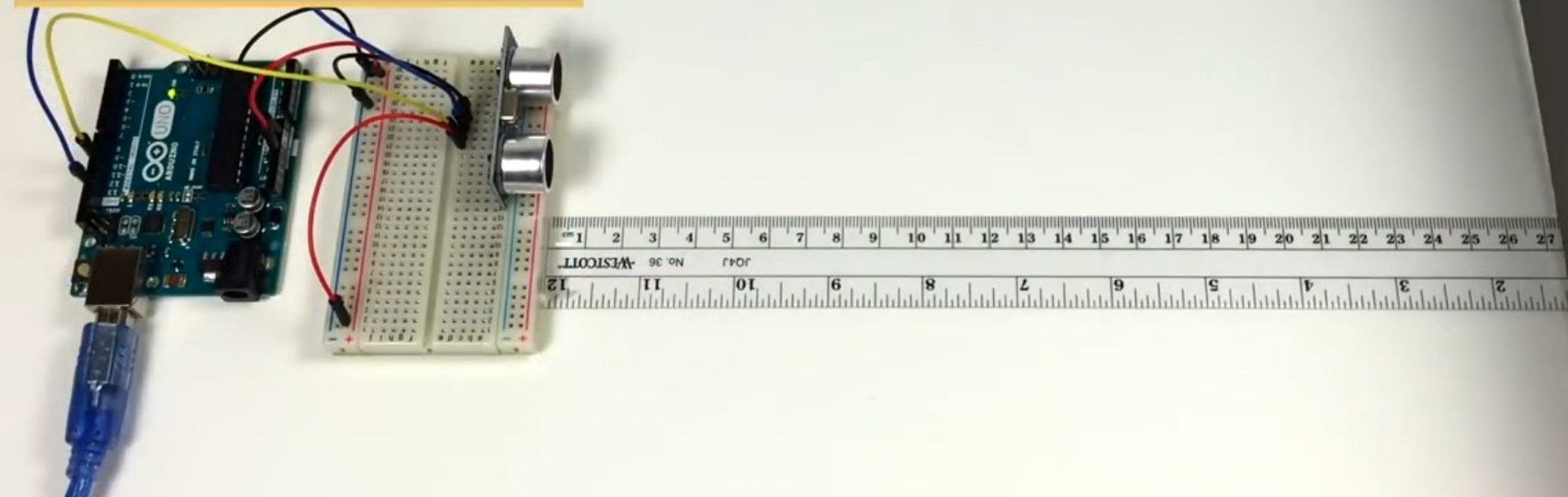
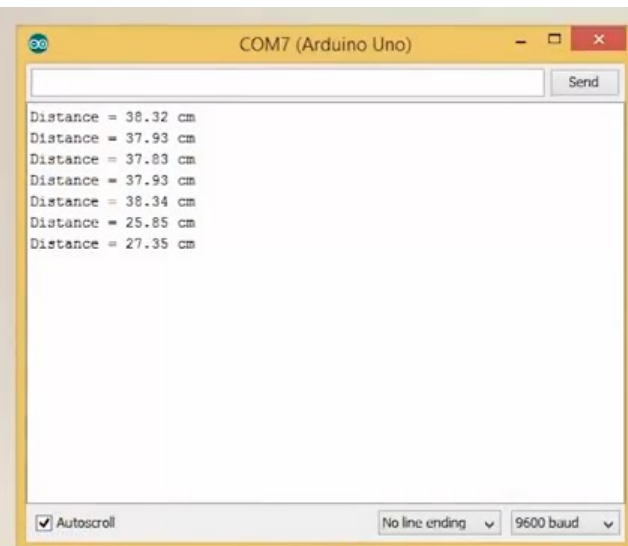


Image credit:
<https://www.parallax.com/product/910-28015a>

Video:
https://youtu.be/Ea4CgAw6b_M?t=683



Our Model

- **(Assumed) Model:**
 - Range reading = true range + error
- Formalize:

$$Y = r + X$$
$$X \sim \text{Norm}_x[\mu, \sigma^2]$$

- **Problem:** Don't know parameters $\theta = \{\mu, \sigma^2\}$
- **Solution:** Learn from data!
 - Fix r to some distance (1m)
 - Collect range reading deviations ($x[i] = y[i] - r$)
 - Estimate (learn) parameters $\theta = \{\mu, \sigma^2\}$

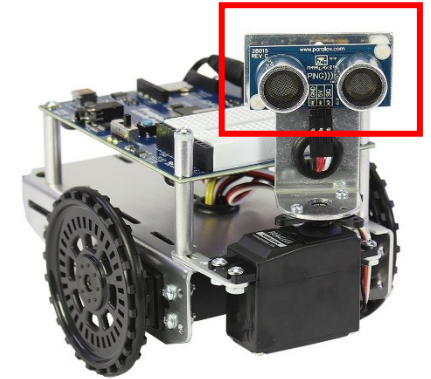
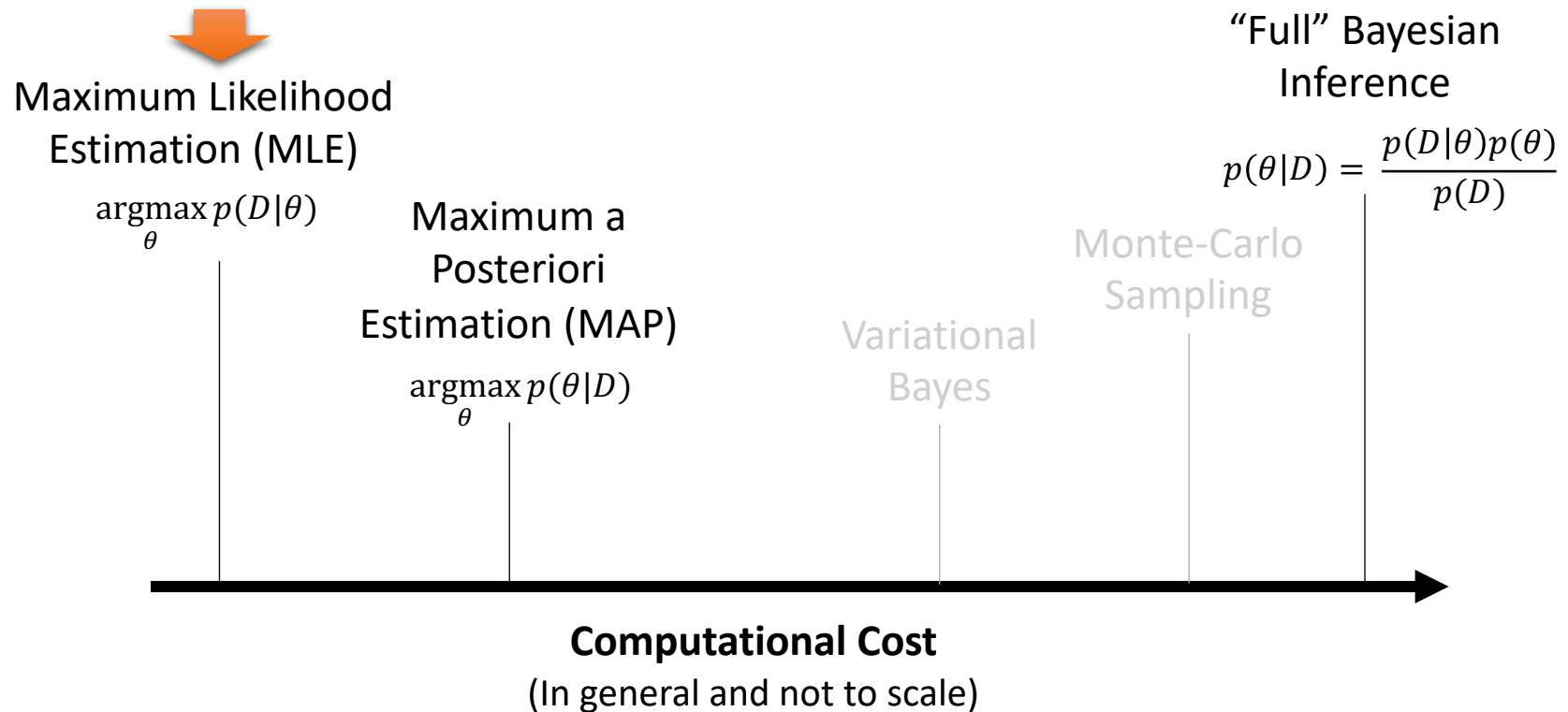


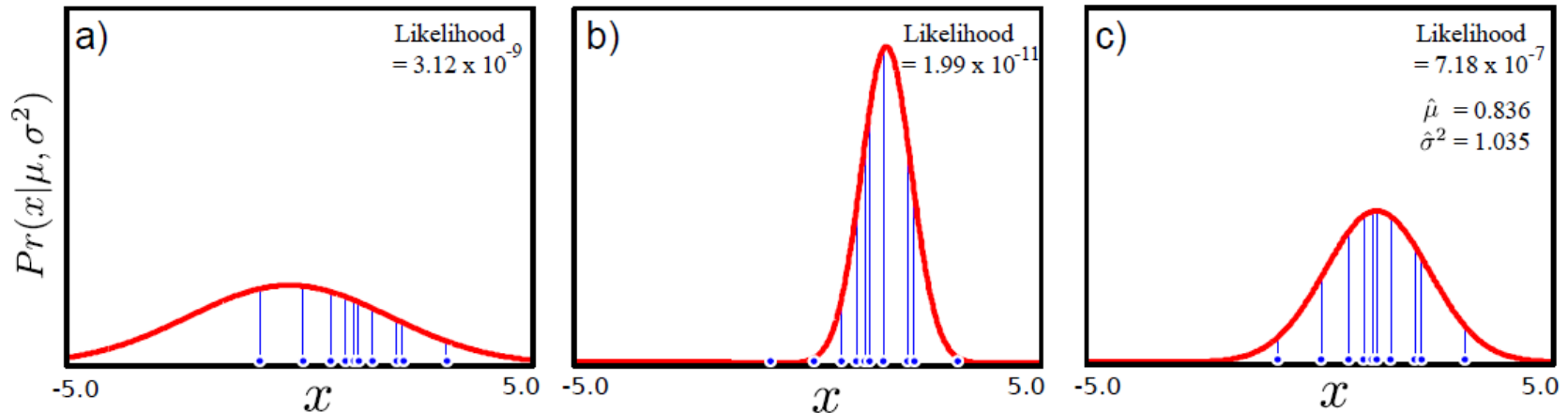
Image credit:
<https://www.parallax.com/product/910-28015a>

Learning Parameters

- Common approaches to **learn the unknown parameters** θ from a set of given data $\mathcal{D} = \{x[1], \dots, x[N]\}$:



Maximum Likelihood Estimate: Intuition

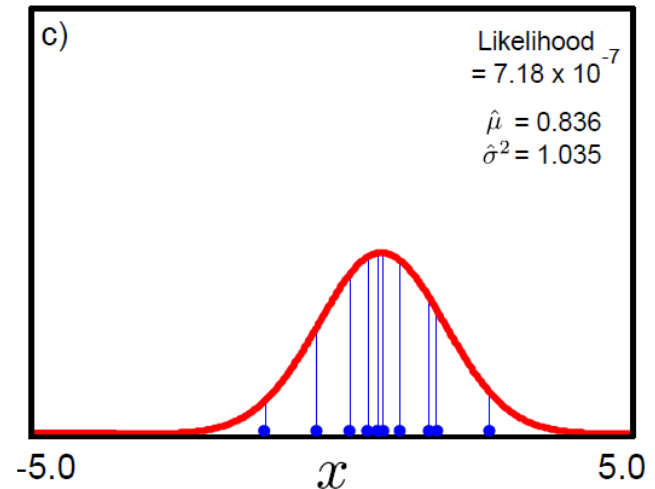
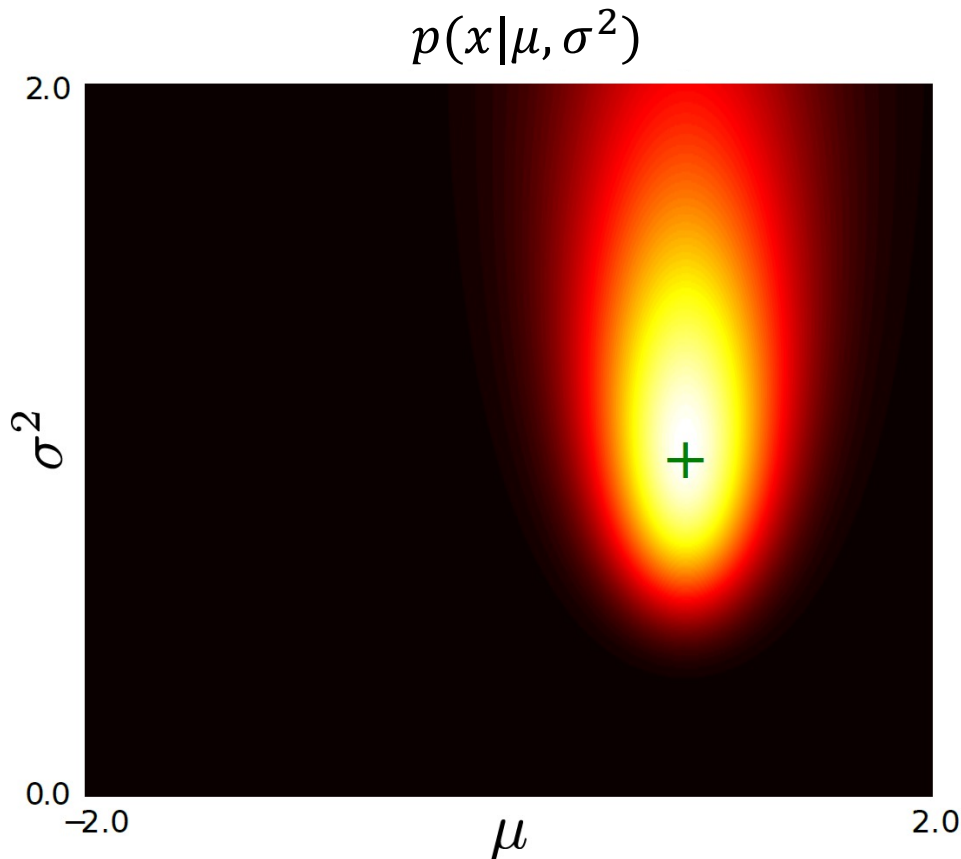


- Blue dots are the observed data $\mathcal{D} = \{x[1], \dots, x[N]\}$.
- Red curves are the Normal distribution for a possible μ and σ^2 .
- The likelihood of a set of **independently** sampled data is the **product** of the individual likelihoods $p(x|\mu, \sigma^2)$ (blue vertical lines).
- The maximum likelihood should be correct μ and σ^2

Example 1: Univariate Normal Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

Intuition behind MLE:



Plotted surface of likelihoods as a function of possible parameter values.

ML Solution is at the **peak**.

Maximum Likelihood Estimation

- Given data $\mathcal{D} = \{x[1], \dots, x[N]\}$
- Assume:
 - a set of distributions $\{p_\theta: \theta \in \Theta\}$ where $p_\theta = p(x|\theta)$
 - \mathcal{D} is sample from $X_1, X_2, \dots, X_N \sim p_{\theta^*}$ for some $\theta^* \in \Theta$
 - Random variables X_1, X_2, \dots, X_N are **independent** and **identically distributed (iid)** according to p_{θ^*}
- Goal: Estimate θ^*
- The estimate θ_{MLE} is a maximum likelihood estimate (MLE) for θ^* if

$$\theta_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} [p(\mathcal{D}|\theta)]$$

Independent and Identically Distributed (iid)

- Common **assumption** in many modeling and learning scenarios
- Allows us to **decompose the likelihood** into products of likelihoods (one for each datum)

$$\begin{aligned}\theta_{MLE} &= \operatorname{argmax}_{\theta} [p(\mathcal{D}|\theta)] \\ &= \operatorname{argmax}_{\theta} [\prod_{i=1}^N p(X = x[i] | \theta)] \quad (\text{i.i.d})\end{aligned}$$

Sensor Uncertainty: MLE

$$\theta_{MLE} = \operatorname{argmax}_{\theta} [\prod_{i=1}^N p(X = x[i] | \theta)]$$

In our case, X is Normal / Gaussian distributed.

Fit an univariate normal distribution model to a set of scalar data $\mathcal{D} = \{x[1], \dots x[N]\}$.

Recall that the univariate normal distribution is given by:

$$p(x) = \operatorname{Norm}_x[\mu, \sigma^2] = \frac{1}{\sqrt{2\pi\sigma^2}} \exp - \frac{(x - \mu)^2}{2\sigma^2}$$

Our goal is to **find the two unknown parameters μ and σ^2** .

Example 1: Univariate Normal Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

$$\begin{aligned}\theta_{MLE} &= \operatorname{argmax}_{\theta} [p(x|\theta)] \\ &= \operatorname{argmax}_{\theta} \left[\prod_{i=1}^N p(x[i] | \theta) \right] \quad (\text{iid})\end{aligned}$$

Likelihood given by pdf

$$p(x|\mu, \sigma^2) = \text{Norm}_x[\mu, \sigma^2] = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x - \mu)^2}{2\sigma^2}$$

Example 1: Univariate Normal Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

Algebraically:

$$\hat{\mu}, \hat{\sigma}^2 = \operatorname{argmax}_{\mu, \sigma^2} [p(x|\mu, \sigma^2)]$$

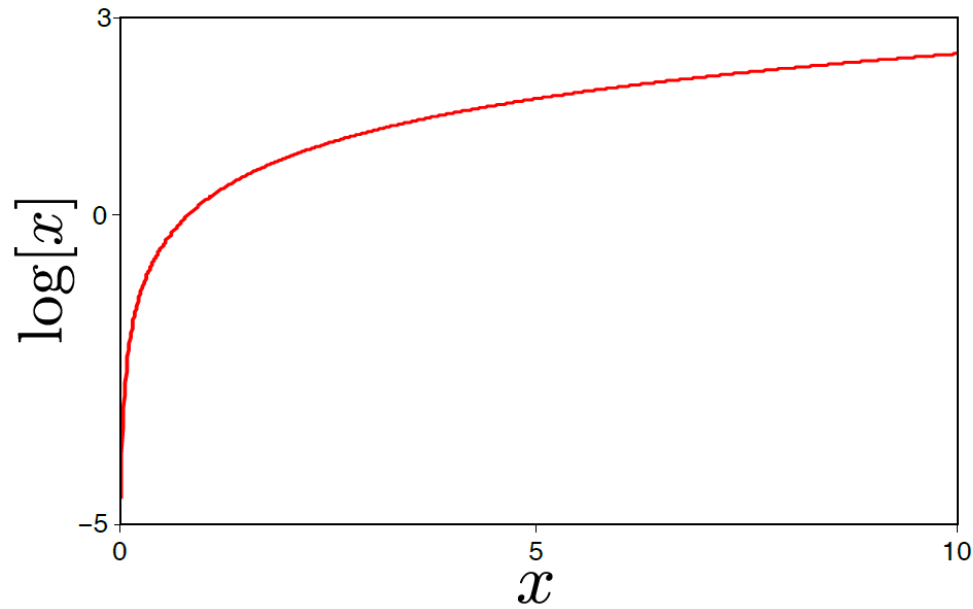
where

$$p(x|\mu, \sigma^2) = \prod_{i=1}^N \operatorname{Norm}_{x[i]} [\mu, \sigma^2],$$

or alternatively, we can maximize the logarithm:

$$\begin{aligned} \hat{\mu}, \hat{\sigma}^2 &= \operatorname{argmax}_{\mu, \sigma^2} \sum_{i=1}^N \log [\operatorname{Norm}_{x[i]} [\mu, \sigma^2]] \\ &= \operatorname{argmax}_{\mu, \sigma^2} \left[-0.5N \log [2\pi] - 0.5N \log \sigma^2 - 0.5 \sum_{i=1}^N \frac{(x[i] - \mu)^2}{\sigma^2} \right] \end{aligned}$$

Why the Logarithm?



- The logarithm is a **monotonic** transformation.
- Hence, the position of the **peak stays in the same place**.
- The log likelihood is **easier to work with**.

Example 1: Univariate Normal Distribution

Approach 1: Maximum Likelihood Estimation (MLE)

$$\begin{aligned}\hat{\mu}, \hat{\sigma}^2 &= \operatorname{argmax}_{\mu, \sigma^2} \sum_{i=1}^N \log [\operatorname{Norm}_{x[i]}[\mu, \sigma^2]] \\ &= \operatorname{argmax}_{\mu, \sigma^2} \underbrace{\left[-0.5N \log [2\pi] - 0.5N \log \sigma^2 - 0.5 \sum_{i=1}^N \frac{(x[i] - \mu)^2}{\sigma^2} \right]}_L\end{aligned}$$

Maximization can be done in closed-form by taking **derivative w.r.t. the variable and equate to zero**:

$$\frac{\partial L}{\partial \mu} = \sum_{i=1}^N \frac{(x[i] - \mu)}{\sigma^2} = \frac{\sum_{i=1}^N x[i]}{\sigma^2} - \frac{N\mu}{\sigma^2} = 0, \quad \frac{\partial L}{\partial \sigma^2} = -\frac{N}{\sigma^2} + \sum_{i=1}^N \frac{(x[i] - \mu)^2}{\sigma^4} = 0$$

$$\Rightarrow \hat{\mu} = \frac{\sum_{i=1}^N x[i]}{N} = \bar{x}, \quad \Rightarrow \hat{\sigma}^2 = \frac{\sum_{i=1}^N (x[i] - \mu)^2}{N}$$

Least Squares Interpretation

Maximum likelihood for the mean of the normal distribution...

$$\hat{\mu} = \operatorname{argmax}_{\mu} \left[-0.5N \log [2\pi] - 0.5N \log \sigma^2 - 0.5 \sum_{i=1}^N \frac{(x[i] - \mu)^2}{\sigma^2} \right]$$

$$= \operatorname{argmax}_{\mu} \left[- \sum_{i=1}^N (x[i] - \mu)^2 \right]$$

$$= \operatorname{argmin}_{\mu} \left[\sum_{i=1}^N (x[i] - \mu)^2 \right]$$

...gives 'least squares' fitting criterion.

Let's try it out.

<https://github.com/crslab/CS5340-notebooks>



MLE: Properties

- **Easy and fast** to compute
- Nice Asymptotic properties:
 - **Consistent**: if data generated from $f(\theta^*)$, MLE converges to its true value, $\hat{\theta}_{MLE} \rightarrow \theta^*$ as $n \rightarrow \infty$
 - **Efficient**: there is no consistent estimator that has lower mean squared error than the MLE estimate (achieves Cramer-Rao lower bound)
- **Functional Invariance**: if $\hat{\theta}$ is the MLE of θ^* , and $g(\theta^*)$ is a transformation of θ^* then the MLE for $\alpha = g(\theta^*)$ is $\hat{\alpha} = g(\hat{\theta})$

What is a problem?

- Imagine if you had samples:
 $\{0.3, -0.1, 1.2, 0.2, -0.2\}$
 - What is your MLE estimate of the mean and variance?
 - $\mu_{MLE} = 0.28, \sigma^2 = 0.245$
- Manual says that on average devices have *zero* bias $\mu = 0$ and variance $\sigma^2 = 0.05$
- How **certain** are you that your estimate is correct?

Other issues

- MLE is a **point estimate** i.e., does not represent uncertainty over the estimate
- MLE **may overfit**.
- MLE **does not incorporate prior information**.
- Asymptotic results are for the **limit** and **assumes model is correct**.
- MLE **may not exist** or **may not be unique**

How can we model our uncertainty about the parameters estimates μ and σ^2 ?



NUS
National University
of Singapore

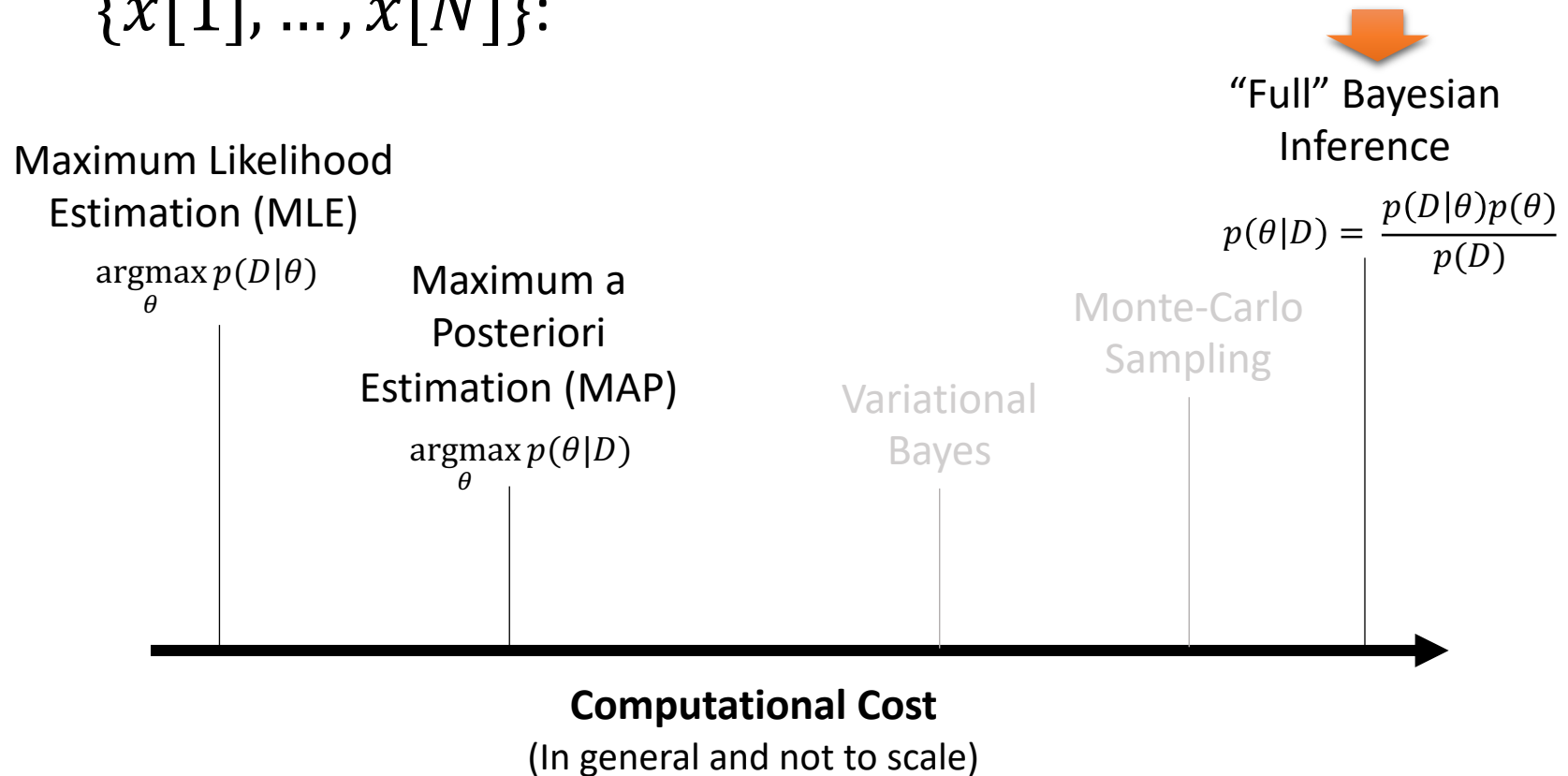
School of
Computing

Learning via Bayes

Bayesian Inference and Conjugate Models

Learning Parameters

- Common approaches to **learn the unknown parameters** θ from a set of given data $\mathcal{D} = \{x[1], \dots, x[N]\}$:



Bayesian Approach

- **Fitting:** Instead of a point estimate $\hat{\theta}$, compute the posterior distribution over **all possible** parameter values using Bayes' rule:

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

- **Principle:** why pick one set of parameters? There are many values that could have explained the data. Try to **capture all of the possibilities**.

Our Model

- Possible (Assumed) Model:
 - Range reading = true range + error
- Formalize:

$$Y = r + X$$
$$X \sim \text{Norm}_x[\mu, \sigma^2]$$

- $\theta = \{\mu, \sigma^2\}$ is now a random variable
- Model **uncertainty** over θ using **prior** distribution(s).
- Then, find posterior:

$$p(\theta|D) = \frac{\prod_{i=1}^N p(x[i] | \theta)p(\theta)}{p(D)}$$

- What can be a **prior distribution**?

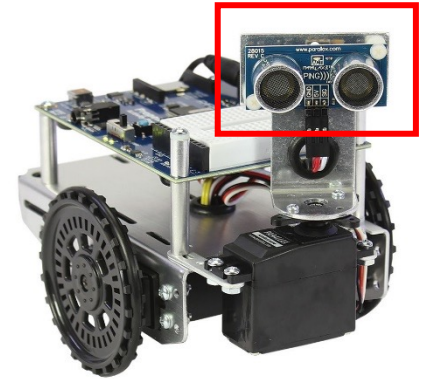


Image credit:
<https://www.parallax.com/product/910-28015a>

Example 1: Univariate Normal Distribution

Approach 3: Bayesian

Compute the posterior distribution using Bayes' rule:

$$p(\theta|x) = \frac{\prod_{i=1}^N p(x[i] | \theta)p(\theta)}{p(x)} = \frac{\prod_{i=1}^N p(x[i] | \theta)p(\theta)}{\int \prod_{i=1}^N p(x[i] | \theta)p(\theta) d\theta}$$

where:

$$\prod_{i=1}^N p(x[i] | \theta)p(\theta) = \prod_{i=1}^N \text{Norm}_{x[i]}[\mu, \sigma^2] \underbrace{\text{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta]}_{\text{Conjugate Prior!}}$$

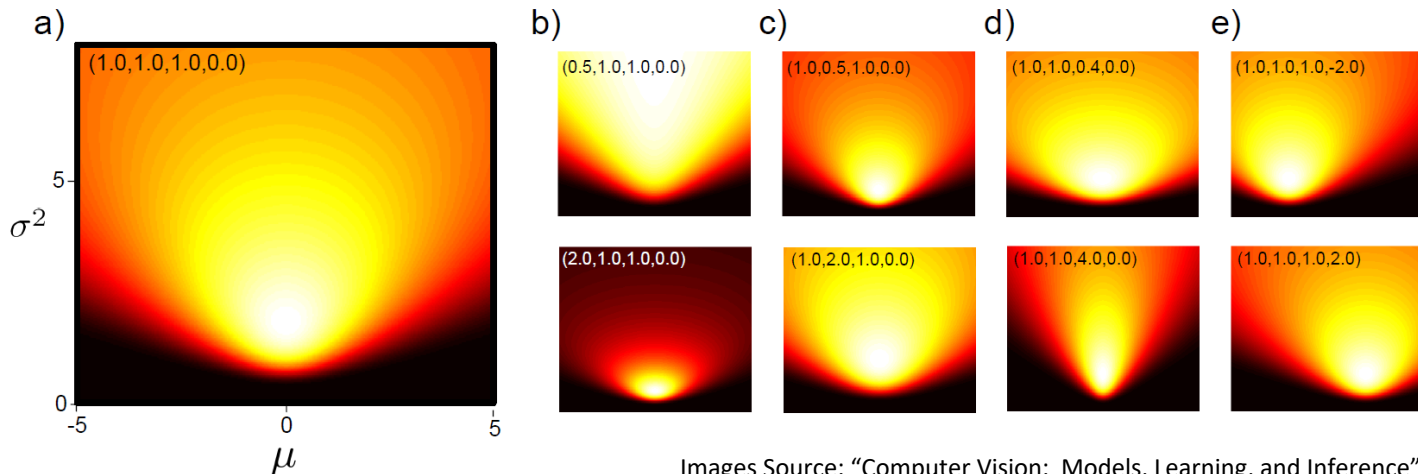
From Lecture 1: Appendix

Normal Inverse Gamma Distribution

$$p(\mu, \sigma^2) = \frac{\sqrt{\gamma}}{\sigma\sqrt{2\pi}} \frac{\beta^\alpha}{\Gamma[\alpha]} \left(\frac{1}{\sigma^2}\right)^{\alpha+1} \exp\left[-\frac{2\beta + \gamma(\delta - \mu)^2}{2\sigma^2}\right]$$

$$p(\mu, \sigma^2) = \text{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta]$$

- **Four hyperparameters** $\alpha, \beta, \gamma > 0$ and $\delta \in \mathbb{R}$.



Images Source: "Computer Vision: Models, Learning, and Inference", Simon Prince

Normal Inverse Gamma

- Where does it “come from”?
 - From Normal and Inverse Gamma!
 - Normal is the prior over mean μ
 - $\text{Norm}_{\mu}[\delta, s]$
 - Inverse Gamma (IG) is the prior over variance σ^2
 - $\text{InvGam}_{\sigma^2}[\alpha, \beta]$
 - Multiply $\text{Norm}_{\mu}[\delta, s]$ and $\text{InvGam}_{\sigma^2}[\alpha, \beta]$ to derive $\text{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta]$
 - where $\gamma = \frac{\sigma^2}{s}$

Example 1: Univariate Normal Distribution

Approach 3: Bayesian

Compute the posterior distribution using Bayes' rule:

$$p(\theta|x) = \frac{\prod_{i=1}^N p(x[i]|\theta)p(\theta)}{p(x)} = \frac{\prod_{i=1}^N p(x[i]|\theta)p(\theta)}{\int \prod_{i=1}^N p(x[i]|\theta)p(\theta) d\theta}$$

$$\prod_{i=1}^N p(x[i]|\theta)p(\theta) = \prod_{i=1}^N \text{Norm}_{x[i]}[\mu, \sigma^2] \text{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta]$$

What distribution is $p(\theta|D)$?

Example 1: Univariate Normal Distribution

Approach 3: Bayesian

Compute the posterior distribution using Bayes' rule:

$$p(\theta|x) = \frac{\prod_{i=1}^N p(x[i]|\theta)p(\theta)}{p(x)} = \frac{\prod_{i=1}^N p(x[i]|\theta)p(\theta)}{\int \prod_{i=1}^N p(x[i]|\theta)p(\theta) d\theta}$$

$$\prod_{i=1}^N p(x[i]|\theta)p(\theta) = \prod_{i=1}^N \text{Norm}_{x[i]}[\mu, \sigma^2] \text{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta]$$

$$p(\theta|D) = \text{NormInvGam}_{\mu, \sigma^2}[\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}, \tilde{\delta}]$$

NormInvGamma is Conjugate Prior for the Normal.

Posterior Form

$$p(\mu, \sigma^2 | D) = \frac{\sqrt{\tilde{\gamma}}}{\sigma \sqrt{2\pi}} \frac{\tilde{\beta}^{\tilde{\alpha}}}{\Gamma[\tilde{\alpha}]} \left(\frac{1}{\sigma^2}\right)^{\tilde{\alpha}+1} \exp\left[-\frac{2\tilde{\beta} + \tilde{\gamma}(\tilde{\delta} - \mu)^2}{2\sigma^2}\right]$$

$$p(\theta | D) = \text{NormInvGam}_{\mu, \sigma^2}[\tilde{\alpha}, \tilde{\beta}, \tilde{\gamma}, \tilde{\delta}]$$

where

$$\tilde{\alpha} = \alpha + \frac{N}{2},$$

$$\tilde{\beta} = \beta + \frac{\sum_i (x[i] - \bar{x})^2}{2} + \frac{N\gamma}{N + \gamma} \frac{(\bar{x} - \delta)^2}{2}.$$

$$\tilde{\delta} = \frac{(\gamma\delta + N\bar{x})}{\gamma + N},$$

$$\tilde{\gamma} = \gamma + N,$$

$$\bar{x} = \frac{1}{N} \sum_i x[i]$$

Let's try it out.



Bayesian Approach: Properties

- Models uncertainty over parameters.
- Principled way of incorporating prior information.
- Can derive quantities of interest, e.g., $p(x < 10|\mathcal{D})$
- Can perform model selection.

Problem

- “Forced” to select a **prior**
- What if your initial belief was **not conjugate** to the normal likelihood?
 - Lognormal, Uniform, Beta ...
- Can be **computationally intractable**

**Can we still incorporate prior information
into the parameter estimation?**

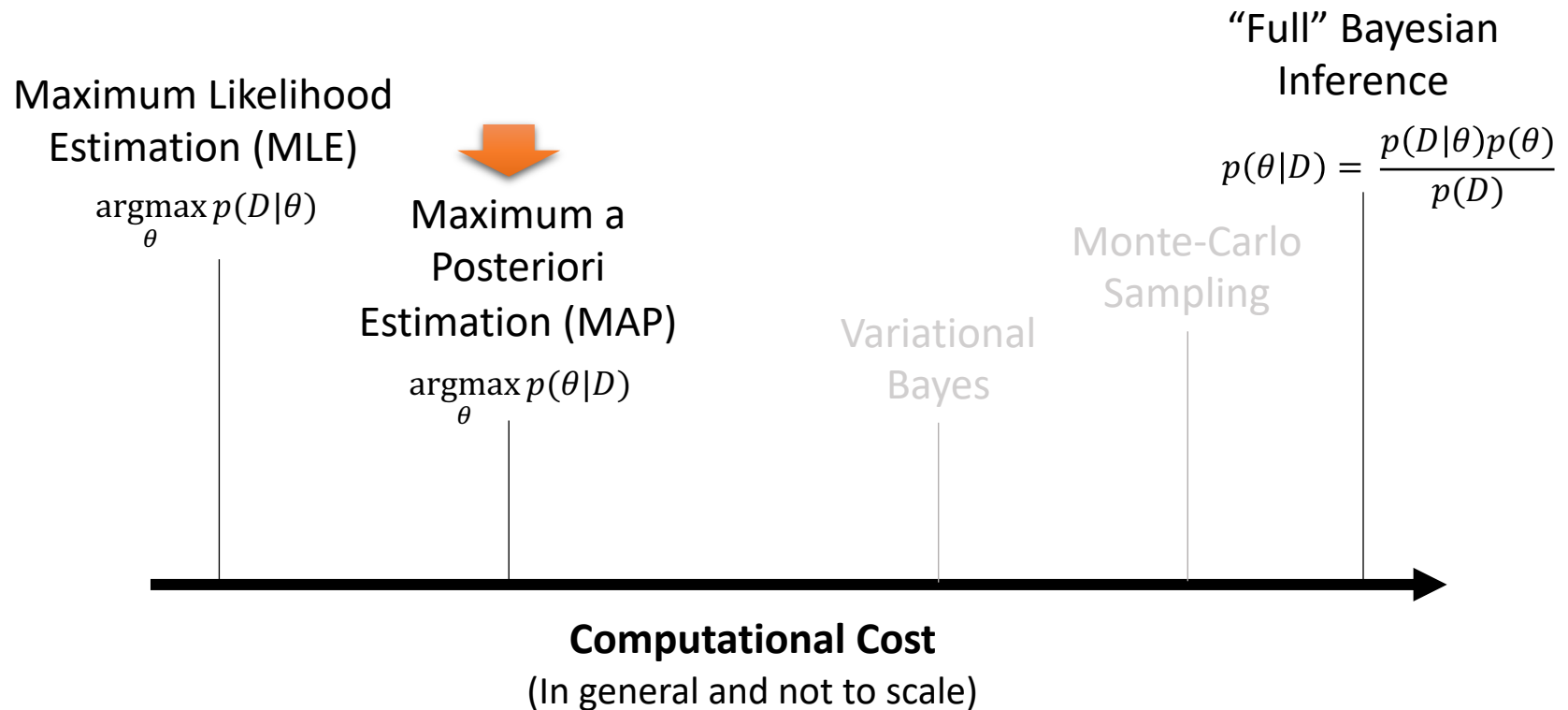
(later in the semester, we will study *approximate* Bayesian inference where we derive an *approximate* posterior distribution)

Learning via MAP

Maximum a Posteriori Estimation(MAP)

Learning Parameters

- Common approaches to **learn the unknown parameters** θ from a set of given data $\mathcal{D} = \{x[1], \dots, x[N]\}$:

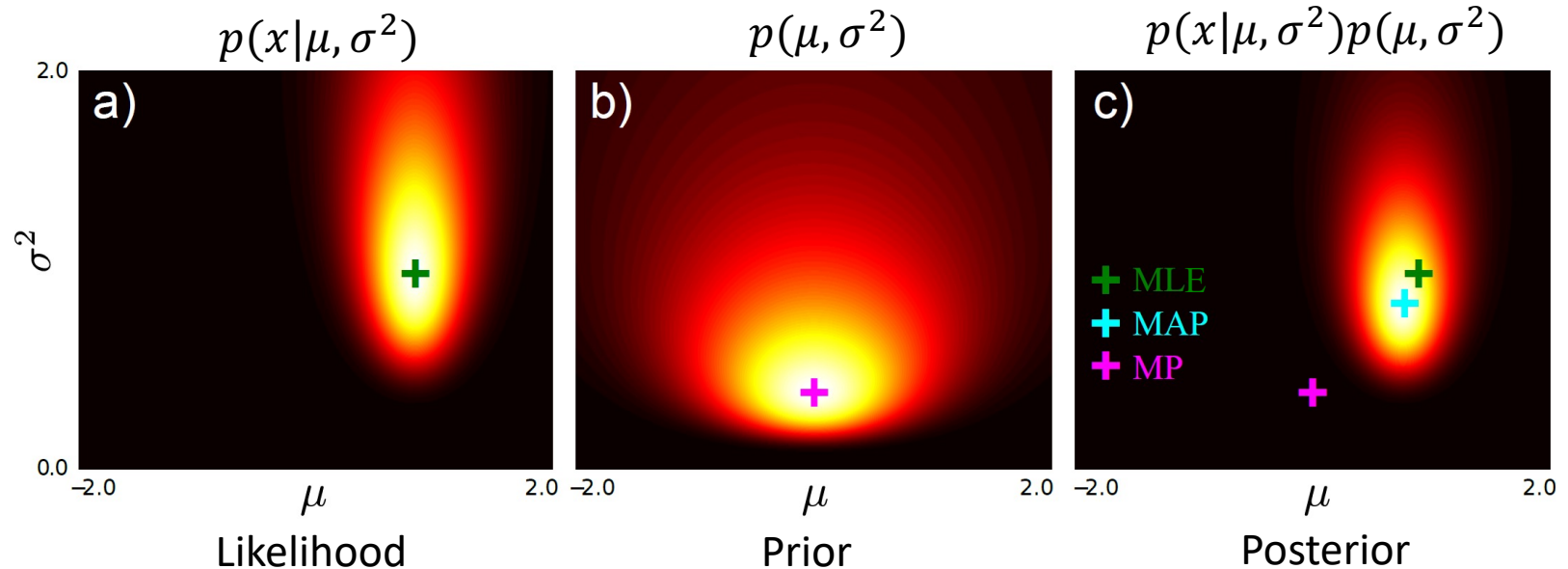


Maximum a Posteriori (MAP)

- Given data $\mathcal{D} = \{x[1], \dots, x[N]\}$
- Assume:
 - Joint distribution $p(\mathcal{D}, \theta)$
 - Here θ is a **random variable**
- Goal: Choose “good” θ
- The estimate θ_{MAP} is a maximum a posteriori estimate (MAP) if

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}}[p(\theta|\mathcal{D})]$$

Intuition: The “Peak” of the Posterior



Maximum a Posteriori (MAP)

- As the name suggests, we find the unknown parameters θ that **maximize the posterior probability** $p(\theta|D)$.

$$\theta_{MAP} = \operatorname{argmax}_{\theta} [p(\theta|D)]$$

$$= \operatorname{argmax}_{\theta} \left[\frac{p(D|\theta)p(\theta)}{p(D)} \right] \quad (\text{Bayes' rule})$$

$$= \operatorname{argmax}_{\theta} \left[\frac{\prod_{i=1}^N p(x[i] | \theta) p(\theta)}{p(D)} \right] \quad (\text{i.i.d})$$

$$= \operatorname{argmax}_{\theta} [\prod_{i=1}^N p(x[i] | \theta) p(\theta)] \quad (p(D) \text{ is removed since it is independent of } \theta)$$

Example 1: Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)

$$\theta_{MAP} = \operatorname{argmax}_{\theta} \left[\prod_{i=1}^N \underbrace{p(x[i] | \theta)}_{\text{Likelihood}} \underbrace{p(\theta)}_{\text{Prior}} \right]$$

Likelihood: univariate Normal distribution

$$p(x|\mu, \sigma^2) = \prod_{i=1}^N \text{Norm}_{x[i]} [\mu, \sigma^2],$$

Prior: normal inverse gamma distribution

$$p(\mu, \sigma^2) = \text{NormInvGam}_{\mu, \sigma^2} [\alpha, \beta, \gamma, \delta]$$

(you can try an alternative prior, but we'll use this for now to compare against the full Bayesian approach)

Example 1: Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)

$$\theta_{MAP} = \underset{\theta}{\operatorname{argmax}} \left[\prod_{i=1}^N \underbrace{p(x[i] | \theta)}_{\text{Likelihood}} \underbrace{p(\theta)}_{\text{Prior}} \right]$$

Likelihood: univariate Normal distribution

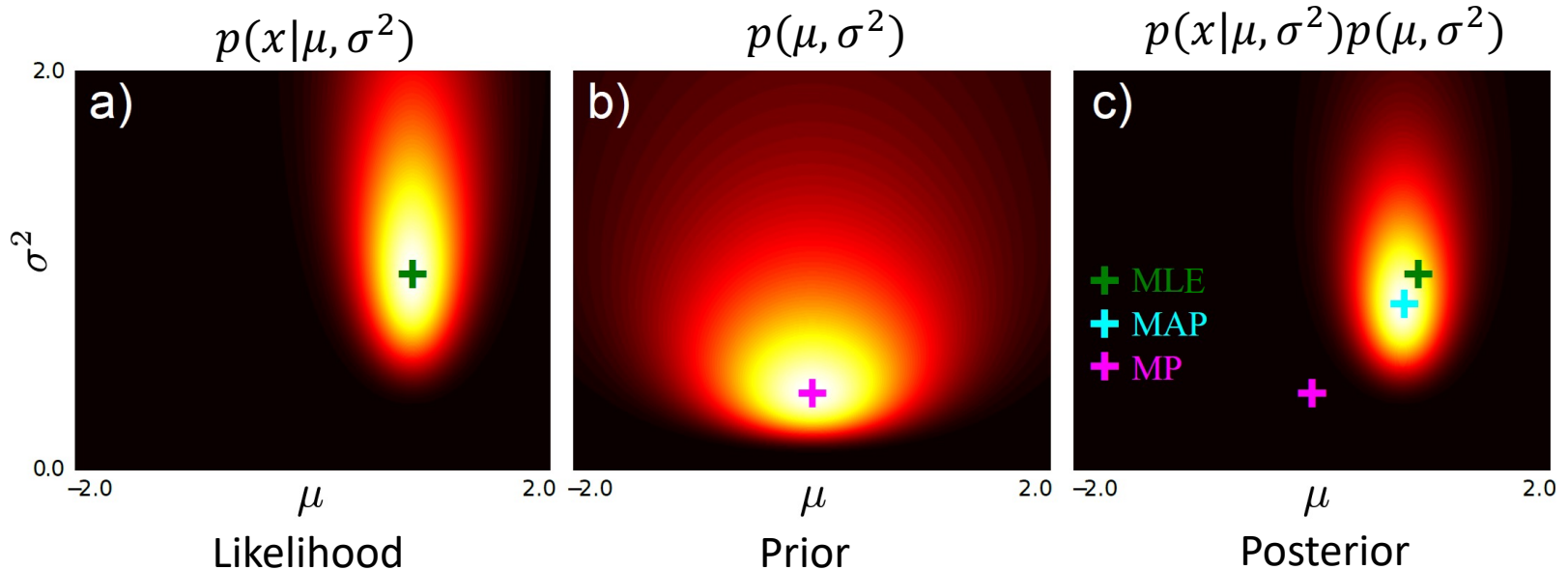
$$p(x|\mu, \sigma^2) = \prod_{i=1}^N \text{Norm}_{x[i]} [\mu, \sigma^2],$$

Prior: normal inverse gamma distribution

$$\begin{aligned} p(\mu, \sigma^2) &= \text{NormInvGam}_{\mu, \sigma^2} [\alpha, \beta, \gamma, \delta] \\ &= \frac{\sqrt{\gamma}}{\sigma \sqrt{2\pi}} \frac{\beta^\alpha}{\Gamma[\alpha]} \left(\frac{1}{\sigma^2} \right)^{\alpha+1} \exp \left[-\frac{2\beta + \gamma(\delta - \mu)^2}{2\sigma^2} \right] \end{aligned}$$

Example 1: Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)



$$\begin{aligned}\hat{\mu}, \hat{\sigma}^2 &= \operatorname{argmax}_{\mu, \sigma^2} \left[\prod_{i=1}^N p(x[i]|\mu, \sigma^2) p(\mu, \sigma^2) \right] \\ &= \operatorname{argmax}_{\mu, \sigma^2} \left[\prod_{i=1}^N \operatorname{Norm}_{x[i]}[\mu, \sigma^2] \operatorname{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta] \right]\end{aligned}$$

Example 1: Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)

$$\begin{aligned}\hat{\mu}, \hat{\sigma}^2 &= \operatorname{argmax}_{\mu, \sigma^2} \left[\prod_{i=1}^N p(x[i] | \mu, \sigma^2) p(\mu, \sigma^2) \right] \\ &= \operatorname{argmax}_{\mu, \sigma^2} \left[\prod_{i=1}^N \operatorname{Norm}_{x[i]}[\mu, \sigma^2] \operatorname{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta] \right]\end{aligned}$$

Maximize the logarithm:

$$\hat{\mu}, \hat{\sigma}^2 = \operatorname{argmax}_{\mu, \sigma^2} \left[\sum_{i=1}^N \log [\operatorname{Norm}_{x[i]}[\mu, \sigma^2]] + \log [\operatorname{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta]] \right]$$

Example 1: Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)

$$\hat{\mu}, \hat{\sigma}^2 = \underset{\mu, \sigma^2}{\operatorname{argmax}} \left[\underbrace{\sum_{i=1}^N \log \left[\operatorname{Norm}_{x[i]}[\mu, \sigma^2] \right]}_L + \log \left[\operatorname{NormInvGam}_{\mu, \sigma^2}[\alpha, \beta, \gamma, \delta] \right] \right]$$

Taking derivatives and setting to zero:

$$\frac{\partial L}{\partial \mu} = 0, \quad \frac{\partial L}{\partial \sigma^2} = 0$$

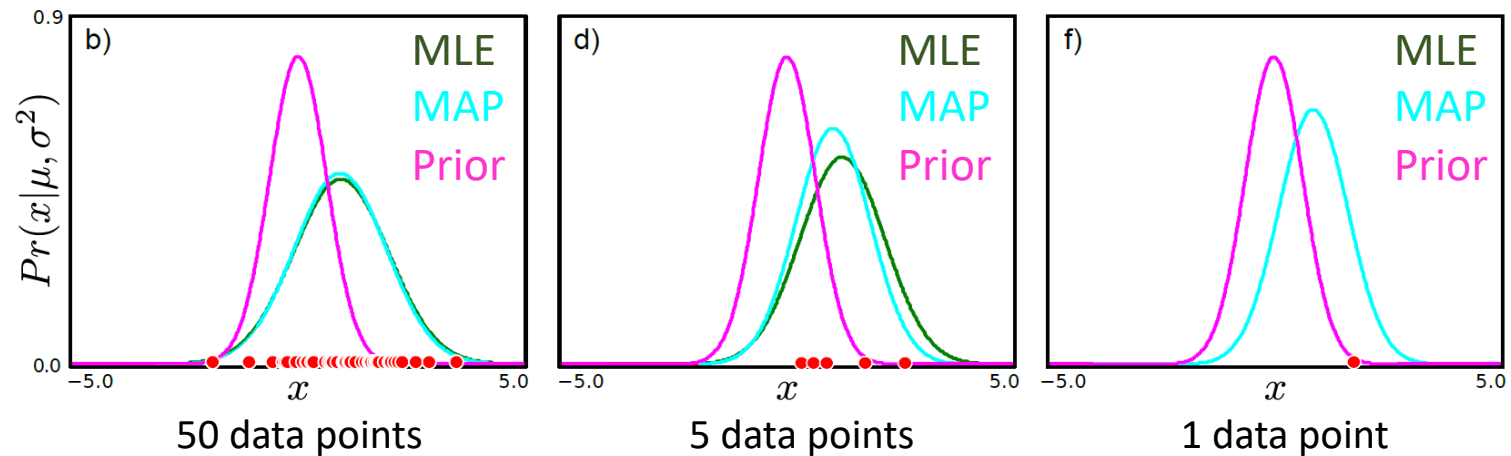
We get:

$$\begin{aligned} \hat{\mu} &= \frac{\sum_i x[i] + \gamma\delta}{N + \gamma}, & \hat{\sigma}^2 &= \frac{\sum_i (x[i] - \hat{\mu})^2 + 2\beta + \gamma(\delta - \hat{\mu})^2}{N + 3 + 2\alpha} \\ &= \frac{N\bar{x} + \gamma\delta}{N + \gamma} \end{aligned}$$

Example 1: Univariate Normal Distribution

Approach 2: Maximum a Posteriori (MAP)

More data points \rightarrow MAP is closer to MLE
Fewer data points \rightarrow MAP is closer to Prior

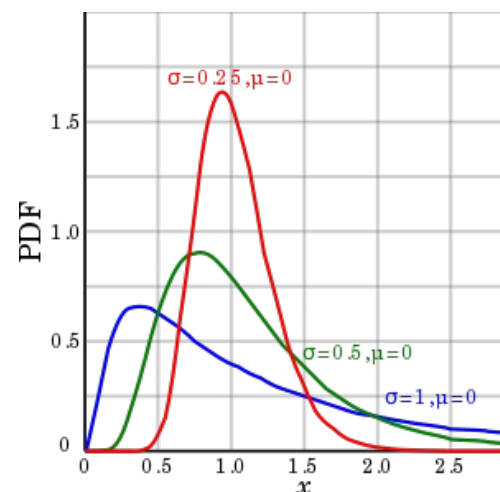


Let's try it out



Take-Home Exercise: Lognormal Prior

- Say you wanted to use a
 - normal prior for μ
 - lognormal prior for σ^2
- Derive the MAP estimates
 - You can derive a closed-form solution for the mean
 - But would need to optimize for σ^2
- Key steps:
 - Derive $\mathcal{L} = \log p(D|\theta)p(\theta) = \log p(D|\theta) + \log p(\theta)$
 - Then set $\frac{\partial \mathcal{L}}{\partial \theta_i} = 0$ for each parameter θ_i



MAP: Properties

- Easy and fast to compute
- Incorporate prior information
- Avoid overfitting (“Regularization”)
- As $n \rightarrow \infty$, MAP tends to look like MLE
 - but does not have the same nice asymptotic properties.

MAP: Problems

- **Point estimate** (like MLE)
 - Does not capture uncertainty over estimates
 - “Poor man’s Bayes”
- Still “forced” to choose prior.
- **NOT Functionally Invariant**: if $\hat{\theta}$ is the MAP of θ^* , and $g(\theta^*)$ is a transformation of θ^* then the MAP for $\alpha = g(\theta^*)$ is **not necessarily** $\hat{\alpha} = g(\hat{\theta})$

Prediction

Maximum Likelihood Estimation (MLE), Maximum a posteriori (MAP), and Bayesian posterior

Predictions for 3 Approaches

Maximum Likelihood Estimate (MLE):

Evaluate new data point x^* under probability distribution with MLE parameters $p(x^*|\theta_{MLE})$.

Maximum a Posteriori (MAP):

Evaluate new data point x^* under probability distribution with MAP parameters $p(x^*|\theta_{MAP})$.

Let's try it out



Predictions for 3 Approaches

Maximum Likelihood Estimate (MLE):

Evaluate new data point x^* under probability distribution with MLE parameters $p(x^*|\theta_{MLE})$.

Maximum a Posteriori (MAP):

Evaluate new data point x^* under probability distribution with MAP parameters $p(x^*|\theta_{MAP})$.

Bayesian:

Calculate weighted sum of predictions from all possible values of parameters

$$p(x^*|D) = \int p(x^*|\theta)p(\theta|D)d\theta$$

Bayesian Approach

Predictive Density:

$$p(x^*|\mathcal{D}) = \frac{p(x^*, D)}{p(D)} \quad \text{(Conditional probability)}$$

$$= \frac{\int p(x^*, D, \theta) d\theta}{p(D)} \quad \text{(Marginalization)}$$

$$= \frac{\int p(x^*, \theta|D) \cancel{p(D)} d\theta}{\cancel{p(D)}} \quad \text{(Chain Rule)}$$

$$= \int p(x^*|D, \theta) p(\theta|D) d\theta \quad \text{(Chain Rule)}$$

$$= \int p(x^*|\theta) p(\theta|D) d\theta \quad \text{(Conditional Independence)}$$

Bayesian Approach

Predictive Density:

$$p(x^*|D) = \int p(x^*|\theta)p(\theta|D)d\theta$$

Prediction for each possible θ Weights

Make a prediction that is an (infinite) **weighted sum** (integral) of the predictions **for each parameter value**, where weights are the probabilities.

Predictive Densities for 3 Approaches

How to rationalize different forms?

Consider MLE and MAP estimates as probability distributions with zero probability everywhere except at estimate (i.e. **delta functions**):

$$\begin{aligned} p(x^*|x) &= \int p(x^*|\theta)\delta[\theta - \hat{\theta}]d\theta \\ &= p(x^*|\hat{\theta}) \end{aligned}$$

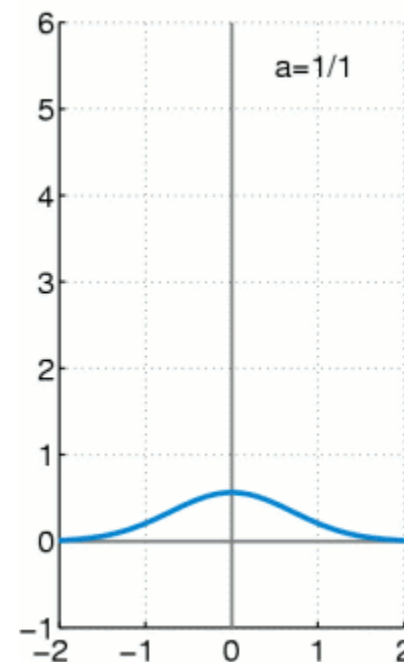
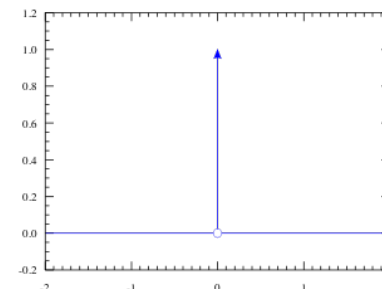


Image Source: Wikipedia

Example 1: Univariate Normal Distribution

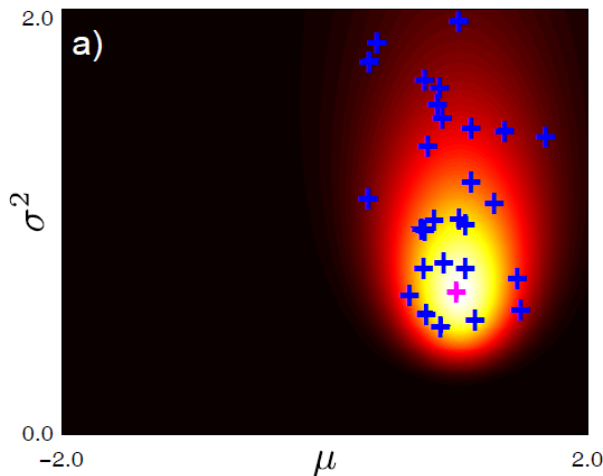
Approach 3: Bayesian

Predictive density

Take weighted sum of predictions from different parameter values:

$$p(x^*|D) = \int \int p(x^*|\mu, \sigma^2) p(\mu, \sigma^2|D) d\mu d\sigma^2$$

Posterior: $p(\mu, \sigma^2|D)$



Samples from posterior

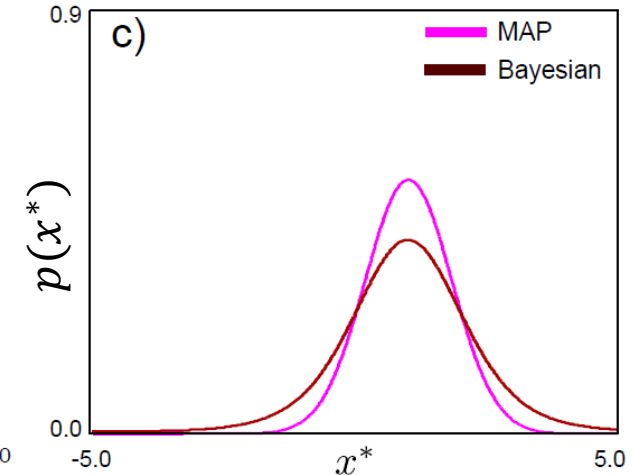
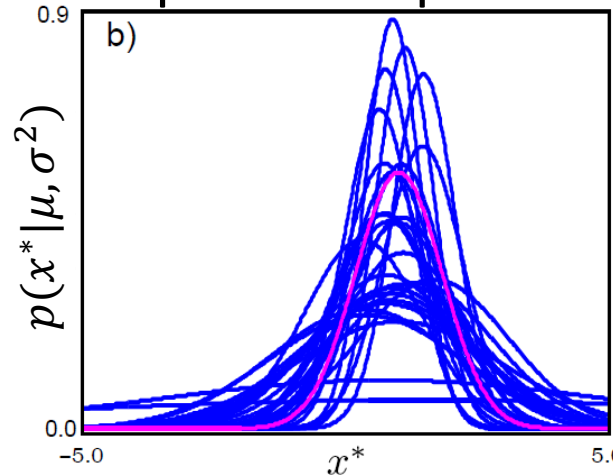


Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince

Example 1: Univariate Normal Distribution

Approach 3: **Bayesian**

Predictive density

Take weighted sum of predictions from different parameter values:

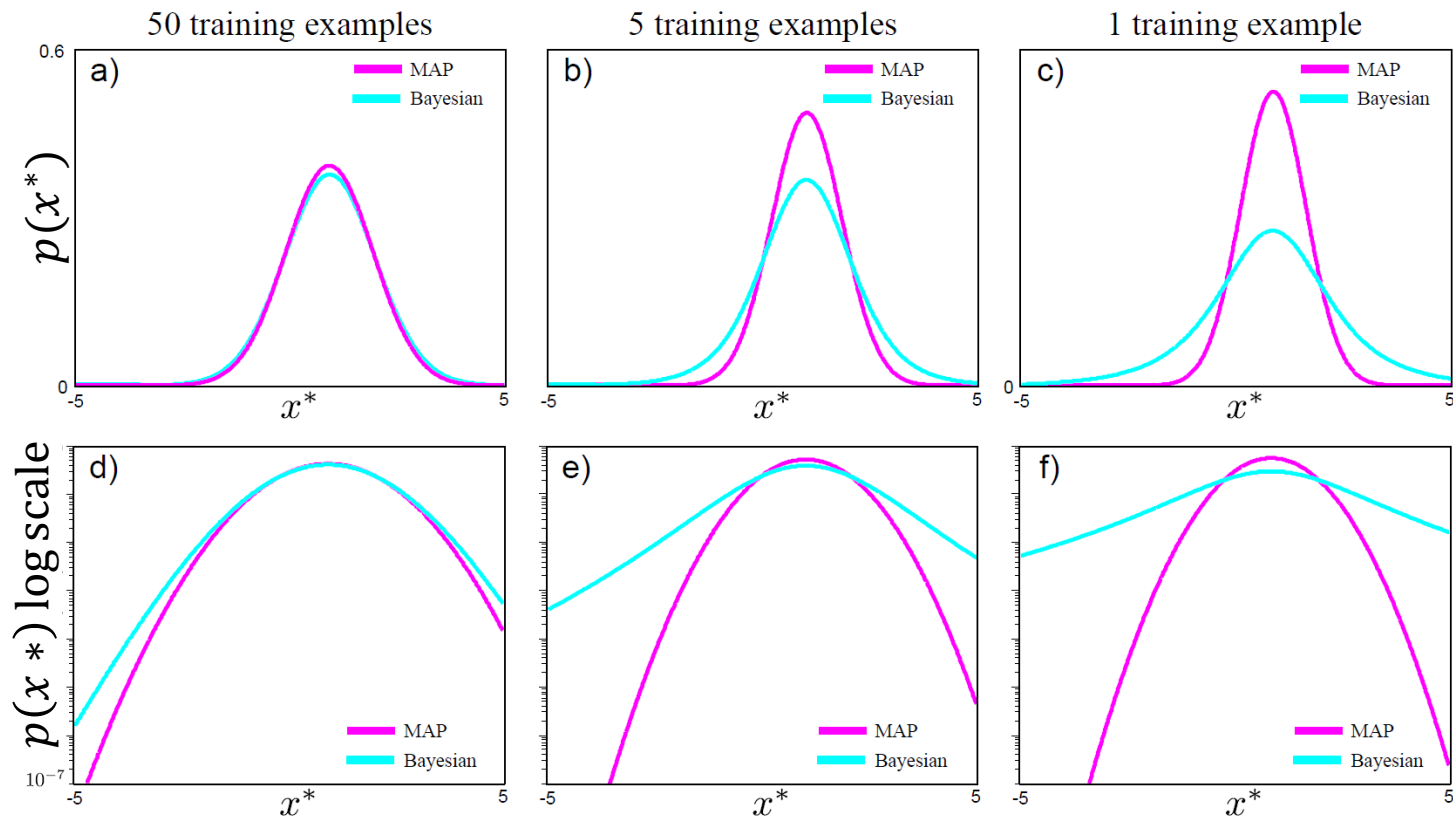
$$p(x^*|x) = t_{2\tilde{\alpha}} \left(x^* | \tilde{\delta}, \frac{\tilde{\beta}(\tilde{\gamma} + 1)}{\tilde{\alpha}\tilde{\gamma}} \right)$$

Where $t_{2\tilde{\alpha}}(x^* | \tilde{\delta}, \frac{\tilde{\beta}(\tilde{\gamma}+1)}{\tilde{\alpha}\tilde{\gamma}})$ is the Generalized Student-T distribution with location $\tilde{\delta}$ and scale $\frac{\tilde{\beta}(\tilde{\gamma}+1)}{\tilde{\alpha}\tilde{\gamma}}$.

Example 1: Univariate Normal Distribution

Approach 3: Bayesian

As the training data decreases, the Bayesian prediction becomes less certain but the MAP prediction can be erroneously overconfident.

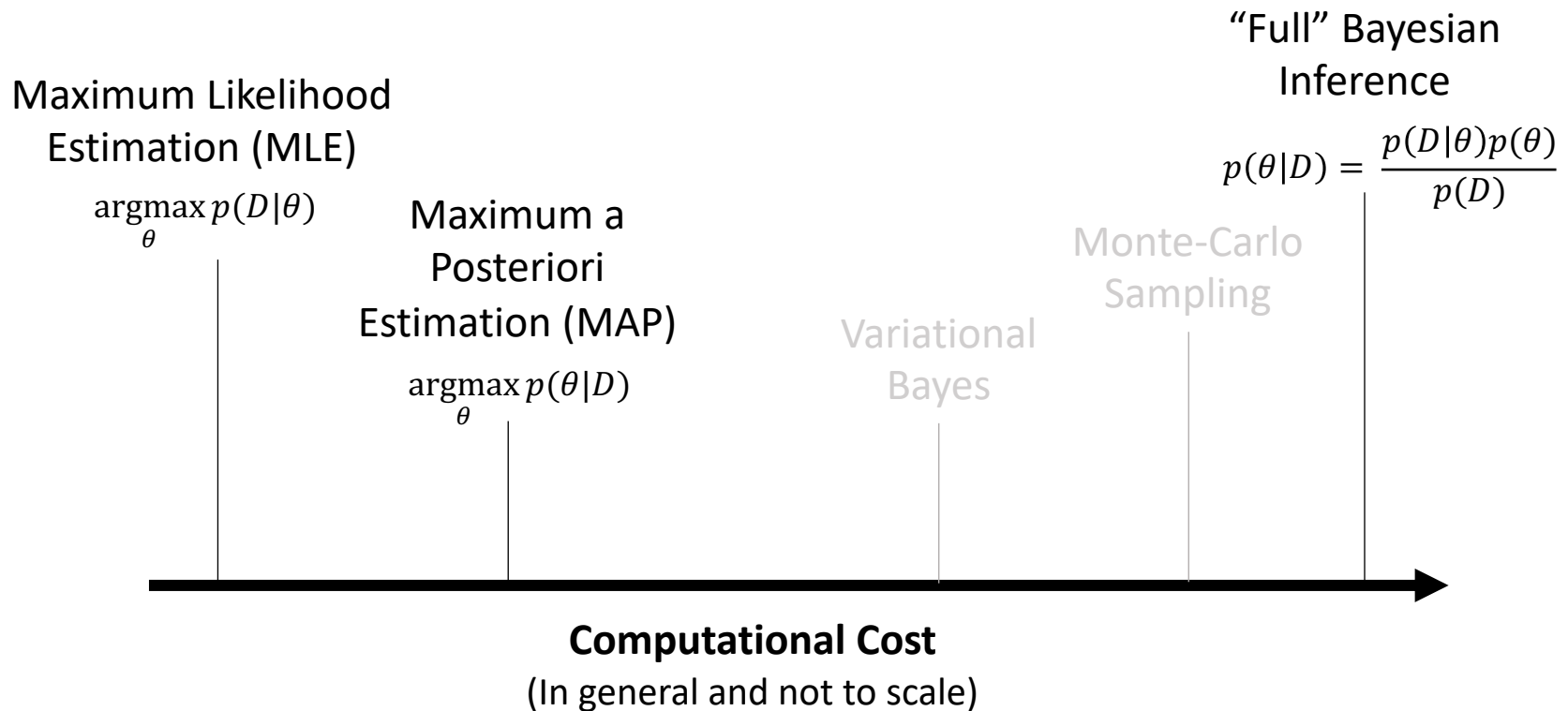


Let's try it out



Learning Parameters

- Common approaches to **learn the unknown parameters** θ from a set of given data $\mathcal{D} = \{x[1], \dots, x[N]\}$:



Exponential Family

What's an Exponential Family and why should we care?

Exponential Family

- An **exponential family** (ExpFam) is a **set** of probability distributions $\{p_\theta: \theta \in \Theta\}$ with the form

$$p_\theta(x) = \frac{h(x) \exp[\eta(\theta)^\top s(x)]}{Z(\theta)}$$

- where:
 - $\theta \in \Theta \subseteq \mathbb{R}^k, x \in \mathbb{R}^d$
 - **Natural parameters:** $\eta(\theta): \Theta \rightarrow \mathbb{R}^m$
 - **Sufficient statistics:** $s(x): \mathbb{R}^d \rightarrow \mathbb{R}^m$
 - **Base Measure (Support and scaling):** $h(x): \mathbb{R}^d \rightarrow [0, \infty)$
 - **Partition function:** $Z(\theta): \Theta \rightarrow [0, \infty)$

Natural/Canonical form

- An exponential family is in its **natural (canonical) form** if it is **parameterized by its natural parameters**:

$$p_{\eta}(x) = p(x|\eta) = \frac{h(x) \exp[\eta^{\top} s(x)]}{Z(\eta)}$$

$$(\text{Compare against } p_{\theta}(x) = \frac{h(x) \exp[\eta(\theta)^{\top} s(x)]}{Z(\theta)})$$

ExpFam: So What?!

- Always **has conjugate prior!**
- Has **fixed number of sufficient statistics** that summarize **iid data** (of **arbitrary** amount!)
- **Posterior predictive distribution** always has **closed form solution** (provided $Z(\theta)$ is closed-form).



The Partition function $Z(\eta)$

$$p_{\eta}(x) = p(x|\eta) = \frac{h(x) \exp[\eta^{\top} s(x)]}{Z(\eta)}$$

- Also called the **normalizer**:

$$Z(\eta) = \int h(x) \exp[\eta^{\top} s(x)] dx$$

- **Why?** To get normalized distribution:

$$\int p(x|\eta) dx = 1$$

$$\int h(x) \frac{\exp[\eta^{\top} s(x)]}{Z(\eta)} dx = 1$$

$$Z(\eta) = \int h(x) \exp[\eta^{\top} s(x)] dx$$

- Aside: sometimes, people write $g(\eta) = 1/Z(\eta)$ and the canonical form becomes:

$$p(x|\eta) = h(x) g(\eta) \exp[\eta^{\top} s(x)]$$

The log Partition function $A(\eta)$

Alternatively, we can specify the log partition function:

$$p_{\eta}(x) = p(x|\eta) = h(x) \exp[\eta^{\top} s(x) - A(\eta)]$$

- Is the **log** of the partition function:

$$A(\eta) = \log Z(\eta) = \log \left[\int h(x) \exp[\eta^{\top} s(x)] dx \right]$$

- **Why?** To get normalized distribution for any η :

$$\int p(x|\eta) dx = \int h(x) \exp[\eta^{\top} s(x) - A(\eta)] = 1$$

$$\exp[-A(\eta)] \int h(x) \exp[\eta^{\top} s(x)] dx = 1$$

$$\exp[A(\eta)] = \int h(x) \exp[\eta^{\top} s(x)] dx$$

$$A(\eta) = \log \left[\int h(x) \exp[\eta^{\top} s(x)] dx \right]$$

Moments of Sufficient Statistics

- For any exponential family distribution:

$$\mathbb{E}[s(x)] = \nabla \log Z(\eta) = \nabla A(\eta)$$

- Higher order moments of $s(x)$ given by higher order derivatives.
- If $s(x) = x$ (**natural exponential family**), we can find moments of x **simply by differentiation!**

MLE of Parameters of ExpFam

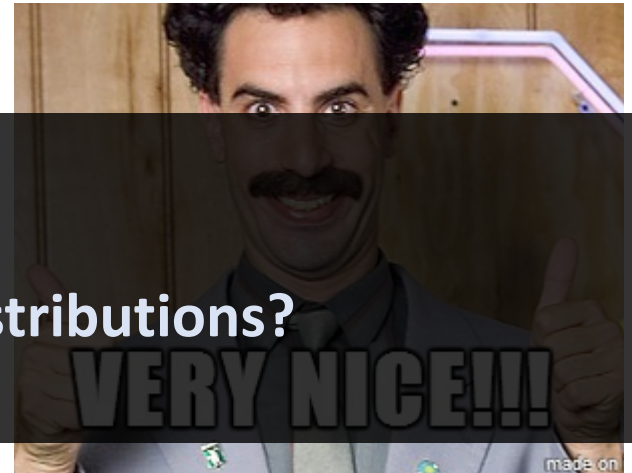
- In addition, the maximum likelihood estimator η_{MLE} satisfies:

$$\nabla A(\eta_{MLE}) = \frac{1}{N} \sum_{n=1}^N s(x_n)$$

- We can use this to solve for η_{MLE}
 - Note that A is convex. Proof in Extra Readings.
- The MLE only depends **only** on **sufficient statistics** $s(x)$

ExpFam: So What?!

- Always **has conjugate prior!**
- Has **fixed number of sufficient statistics** that summarize iid data (of arbitrary amount!) **Great!**
- **Posterior predictive distribution** always has **closed form solution** (provided $Z(\theta)$ is closed-form).



Is Gaussian an ExpFam?

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x - \mu)^2}{2\sigma^2}$$
$$p(x) = \text{Norm}_x[\mu, \sigma^2]$$

Rearrange to fit the ExpFam form:

$$p_\eta(x) = p(x|\eta) = \frac{h(x) \exp[\eta^\top s(x)]}{Z(\eta)}$$

the idea is to match the terms into the:

Natural parameters: $\eta(\theta)$

Sufficient statistics: $s(x)$

Base measure: $h(x)$

Partition function: $Z(\eta)$

or **Log Partition function:** $A(\eta)$

$$p(x|\eta) = h(x) \exp[\eta^\top s(x) - A(\eta)]$$

Many Distributions are ExpFam

PMFs

- Bernoulli
- Binomial
- Categorical/Multinoulli
- Poisson
- Multinomial
- Negative Binomial
- ...

PDFs

- Normal
- Gamma & Inverse Gamma
- Wishart & Inverse Wishart
- Beta
- Dirichlet
- lognormal
- Exponential
- ...

Exercise: Find a family of distributions that is not ExpFam.

Exponential Family

- An **exponential family** (ExpFam) is a **set** of probability distributions $\{p_\theta: \theta \in \Theta\}$ with the form

$$p_\theta(x) = \frac{h(x) \exp[\eta(\theta)^\top s(x)]}{Z(\theta)}$$

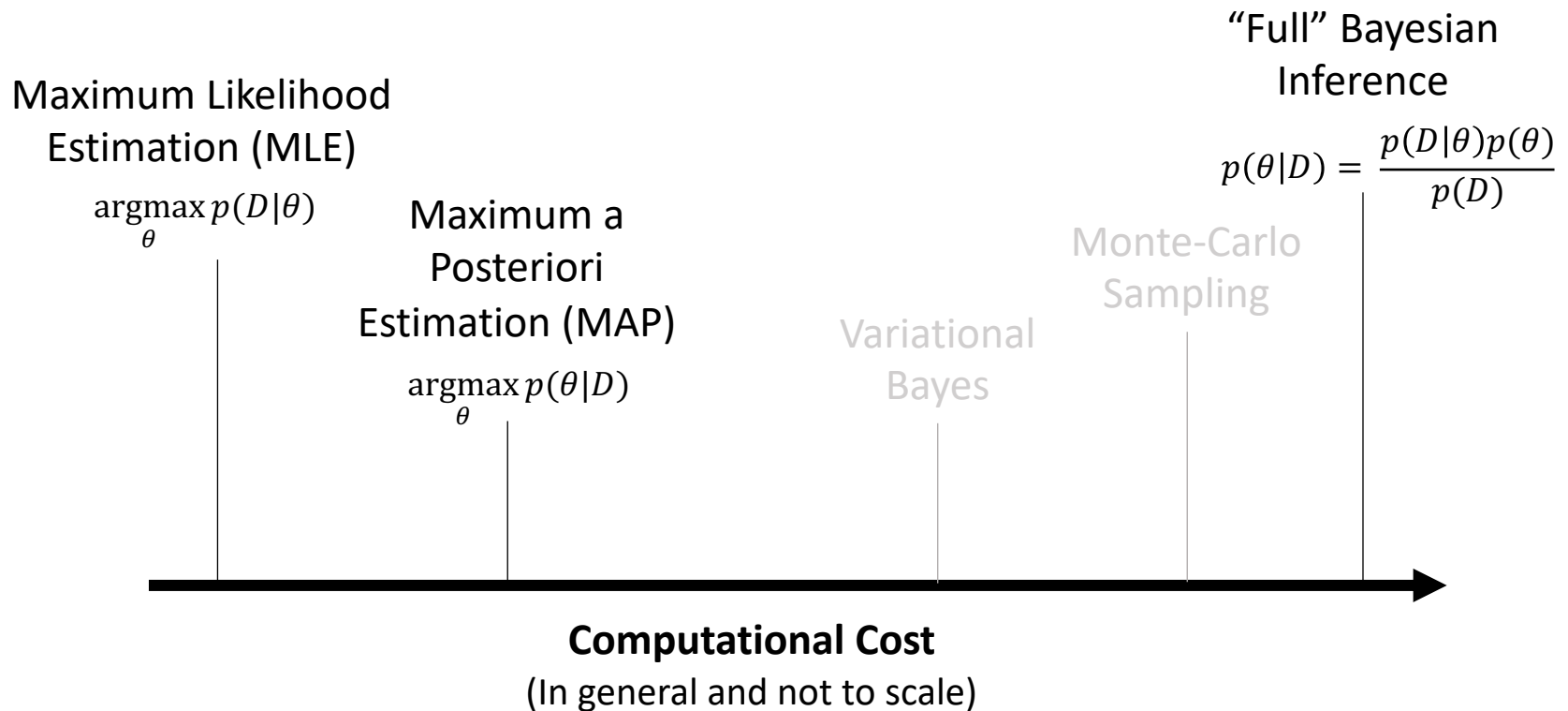
- where:
 - $\theta \in \Theta \subseteq \mathbb{R}^k, x \in \mathbb{R}^d$
 - **Natural parameters:** $\eta(\theta): \Theta \rightarrow \mathbb{R}^m$
 - **Sufficient statistics:** $s(x): \mathbb{R}^d \rightarrow \mathbb{R}^m$
 - **Base Measure (Support and scaling):** $h(x): \mathbb{R}^d \rightarrow [0, \infty)$
 - **Partition function:** $Z(\theta): \Theta \rightarrow [0, \infty)$

Recap

*MLE, MAP, Bayesian Inference, and
Exponential Families*

Learning Parameters

- Common approaches to **learn the unknown parameters** θ from a set of given data $\mathcal{D} = \{x[1], \dots, x[N]\}$:



Exponential Family

- An **exponential family** (ExpFam) is a **set** of probability distributions $\{p_\theta: \theta \in \Theta\}$ with the form

$$p_\theta(x) = \frac{h(x) \exp[\eta(\theta)^\top s(x)]}{Z(\theta)}$$

- where:
 - $\theta \in \Theta \subseteq \mathbb{R}^k, x \in \mathbb{R}^d$
 - **Natural parameters:** $\eta(\theta): \Theta \rightarrow \mathbb{R}^m$
 - **Sufficient statistics:** $s(x): \mathbb{R}^d \rightarrow \mathbb{R}^m$
 - **Base Measure (Support and scaling):** $h(x): \mathbb{R}^d \rightarrow [0, \infty)$
 - **Partition function:** $Z(\theta): \Theta \rightarrow [0, \infty)$

Learning Outcomes

- Students should be able to:
 1. Use the **Maximum Likelihood**, **Maximum a Posteriori** and **Bayesian** approaches to learn the unknown parameters of probability distributions of a **single random variable** from data.
 2. Apply the assumption **independent and identically distributed samples** to simplify the parameter learning process.
 3. Apply the learned parameters to **make predictions**.
 4. Describe the **exponential family** and its properties

A Discrete Example: CS5340 Meme of the Year

CS5340 student: Let me just skip solving tutorials.

screws up in the final exam

CS5340 student:



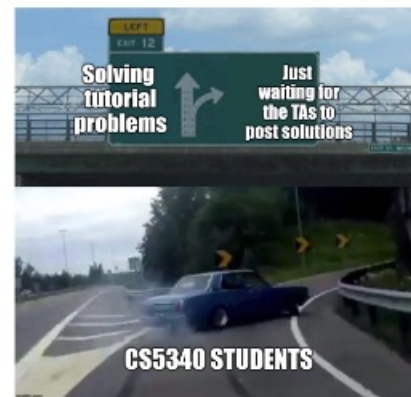
(a) Surprised Pikachu



(b) Two Buttons Dilemma



(c) Distracted Boyfriend



(d) Left Exit 12

CS5340 Meme of the Year

- Model and learn parameters

ID	Template Name	# Votes
1	Surprised Pikachu	25
2	Two Buttons Dilemma	12
3	Distracted Boyfriend	30
4	Left Exit 12	10

Table 1: Votes received by each template by CS5340 students