

#### CS5340 Uncertainty Modeling in Al

Lecture 2: Fitting Probability Models

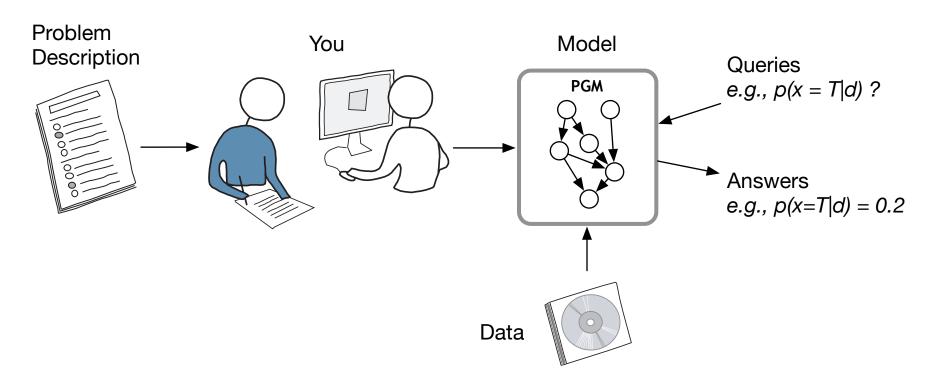
Asst. Prof. Harold Soh
AY 23/24
Semester 2

# Course Schedule (Tentative)

Week	Date	Lecture Topic	Tutorial
1	16 Jan	Introduction to Uncertainty Modeling + Probability Basics	Introduction-
2	23 Jan	Simple Probabilistic Models	Introduction and Probability Basics
3	30 Jan	Bayesian networks (Directed graphical models)	More Basic Probability
4	6 Feb	Markov random Fields (Undirected graphical models)	DGM modelling and d-separation
5	13 Feb	Variable elimination and belief propagation	MRF + Sum/Max Product
6	20 Feb	Factor graphs	Quiz 1
-	-	RECESS WEEK	
7	5 Mar	Mixture Models and Expectation Maximization (EM)	Linear Gaussian Models
8	12 Mar	Hidden Markov Models (HMM)	Probabilistic PCA
9	19 Mar	Monte-Carlo Inference (Sampling)	Linear Gaussian Dynamical Systems
10	26 Mar	Variational Inference	MCMC + Sequential VAE
11	2 Apr	Inference and Decision-Making	Diffusion Models
12	9 Apr	Gaussian Processes (Special Topic)	Quiz 2
13	16 Apr	Project Presentations	Closing Lecture



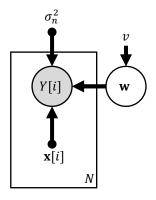
CS5340 is about how to "represent" and "reason" with uncertainty in a computer.





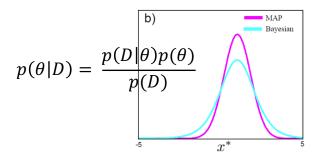
CS5340 :: Harold Soh

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.

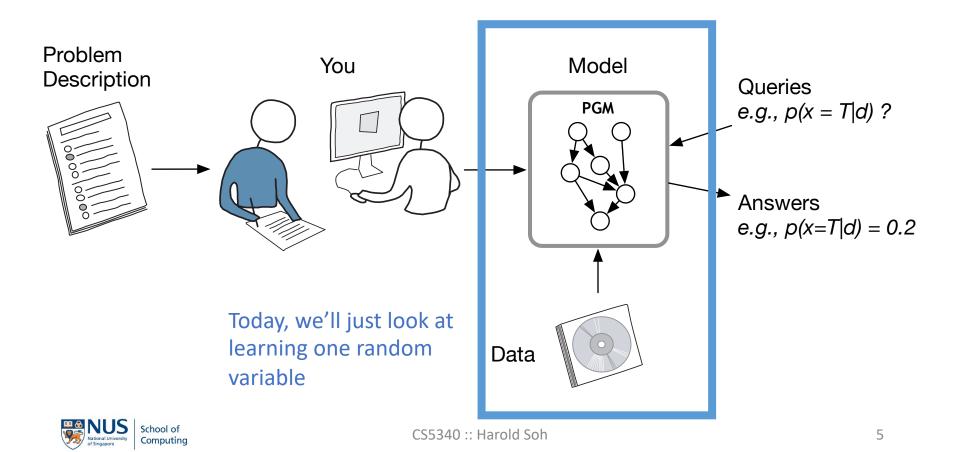


**Reasoning**: We use learning and inference algorithms to answer questions.

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



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:

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

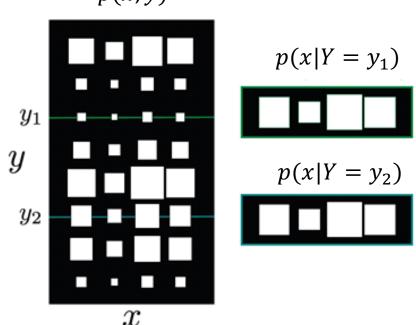
#### Probability: Independence

• The independence of *X* and *Y* means that every conditional distribution is the same.

• The value of Y tells us nothing about X and viceversa. p(x,y)

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

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





# Probability: Bayes' Rule

Recall:

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

• Eliminating p(x, y), we get:

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



Thomas Bayes

• 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

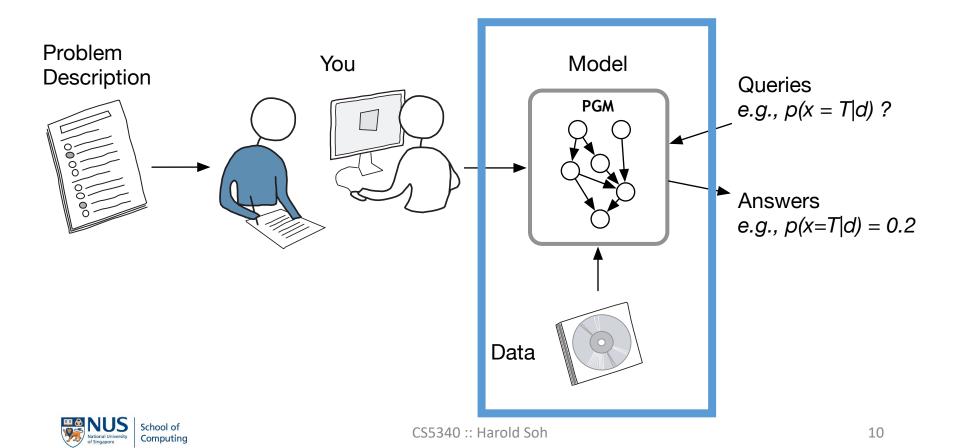
• The expected or average value of some function f[x] taking into account the distribution of X.

#### **Definition:**

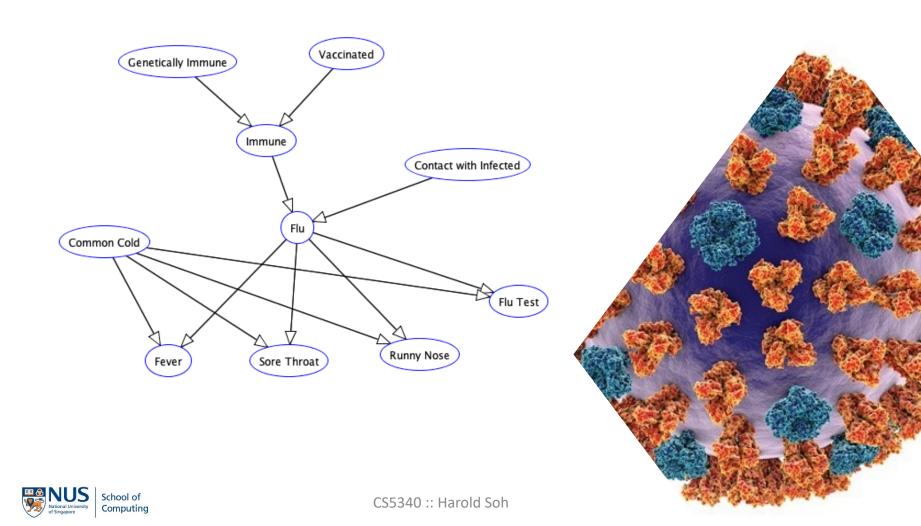
$$E[f[x]] = \sum_{x} f[x]p(x)$$
$$E[f[x]] = \int_{x} f[x]p(x)dx$$



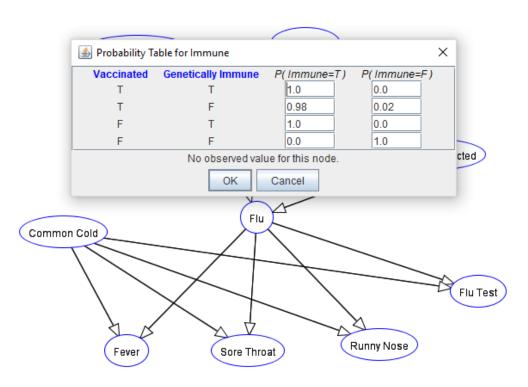
CS5340 is about how to "represent" and "reason" with uncertainty in a computer.

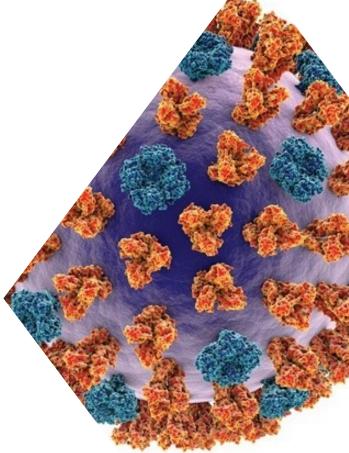


# Generative (Causal) Modeling of Relationships between Variables



# Generative (Causal) Modeling of Relationships between Variables



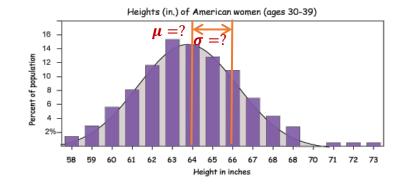




CS5340 :: Harold Soh

#### Fitting Probability Models

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





CS5340 :: Harold Soh 13

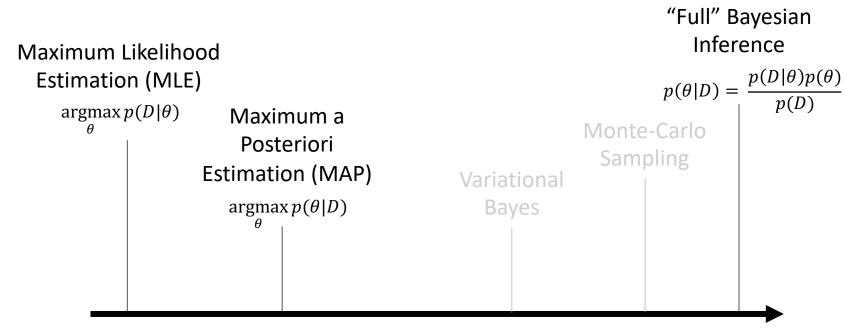
#### Learning Outcomes

- Students should be able to:
  - 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.
  - 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  $\theta$  from a set of given data  $\mathcal{D} = \{x[1], ..., x[N]\}$ :



#### **Computational Cost**

(In general and not to scale)



#### 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.



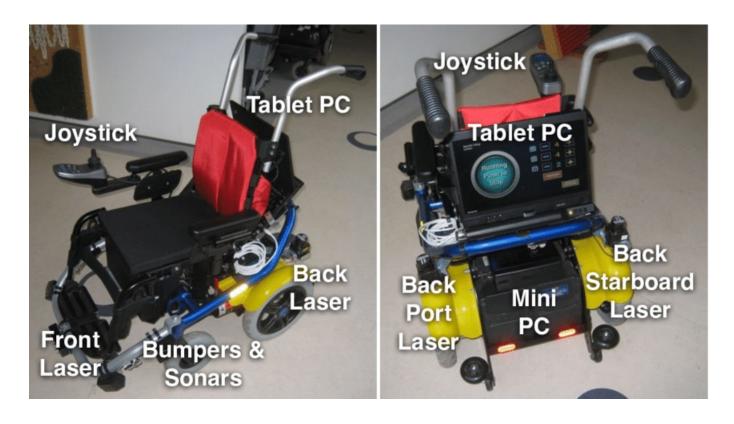


# Learning via MLE

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?

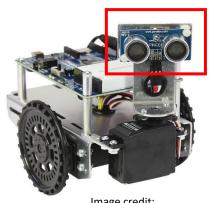
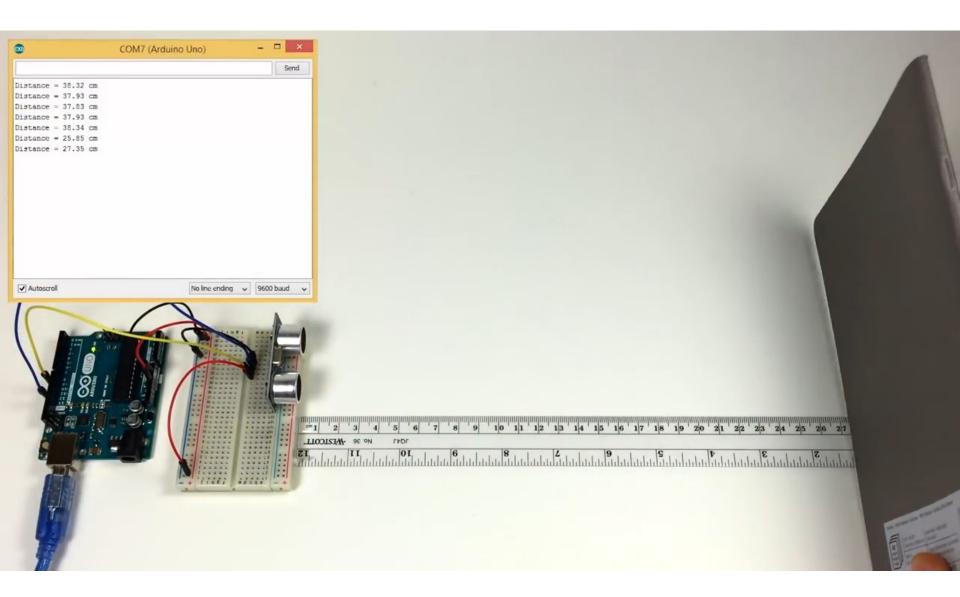


Image credit: https://www.parallax.com/product/910-28015a

Video:

https://youtu.be/Ea4C GAw6b M?t=683







CS5340 :: Harold Soh

#### Our Model

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

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



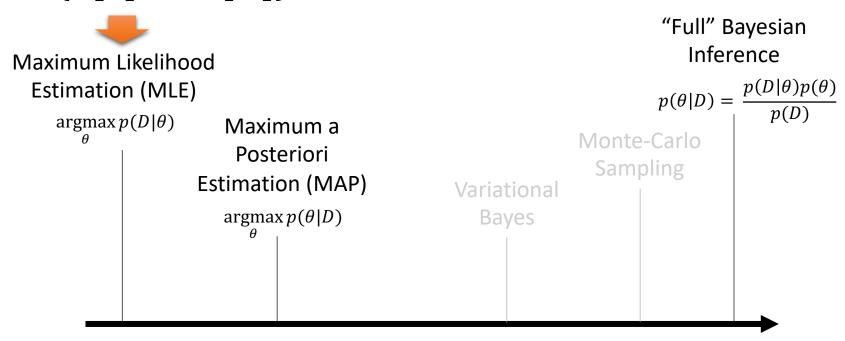
Image credit: https://www.parallax.com /product/910-28015a

- **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  $\theta = \{\mu, \sigma^2\}$



#### Learning Parameters

• Common approaches to learn the unknown parameters  $\theta$  from a set of given data  $\mathcal{D} = \{x[1], ..., x[N]\}$ :

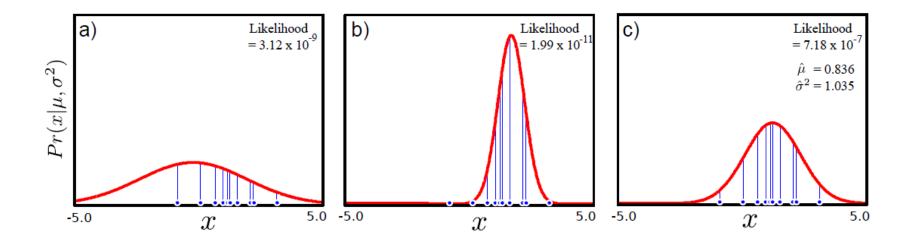


#### **Computational Cost**

(In general and not to scale)



#### Maximum Likelihood Estimate: Intuition



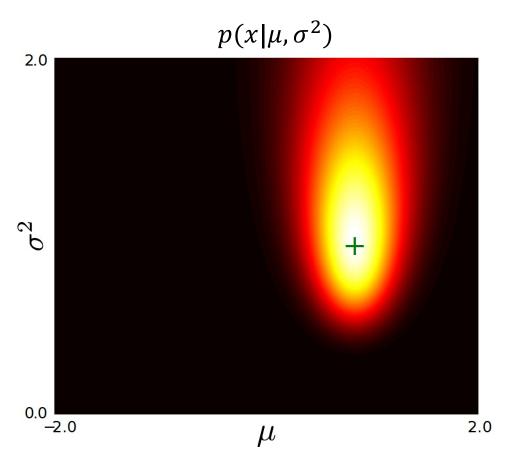
- Blue dots are the observed data  $\mathcal{D} = \{x[1], ..., x[N]\}.$
- Red curves are the Normal distribution for a possible  $\mu$  and  $\sigma^2$ .
- The likelihood of a set of **independently** sampled data is the **product** of the individual likelihoods  $p(x|\mu, \sigma^2)$  (blue vertical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

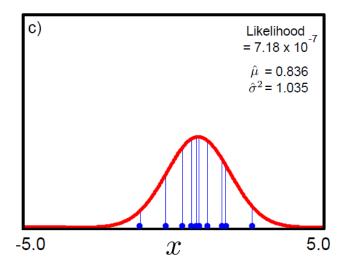


#### Example 1: Univariate Normal Distribution

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

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

$$= \underset{\theta}{\operatorname{argmax}} [\prod_{i=1}^{N} p(X = x[i]|\theta)] \quad \text{(i.i.d)}$$



# Sensor Uncertainty: MLE

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

In our case, X is Normal / Gaussian distributed.

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

Recall that the univariate normal distribution is given by:

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

Our goal is to find the two unknown parameters  $\mu$  and  $\sigma^2$ .



#### Example 1: Univariate Normal Distribution

#### Approach 1: Maximum Likelihood Estimation (MLE)

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} [p(x|\theta)]$$

$$= \underset{\theta}{\operatorname{argmax}} \left[ \prod_{i=1}^{N} p(x[i] \mid \theta) \right]$$
 (iid)

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:

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

where

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

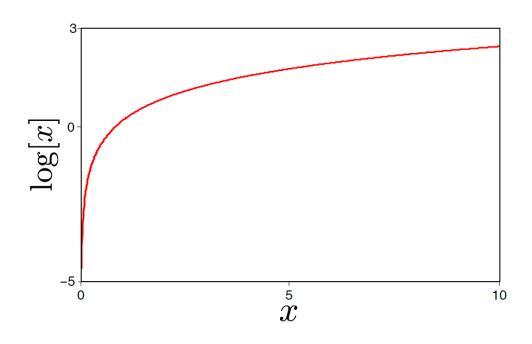
or alternatively, we can maximize the logarithm:

$$\hat{\mu}, \hat{\sigma}^2 = \underset{\mu, \sigma^2}{\operatorname{argmax}} \sum_{i=1}^{N} \log \left[ \operatorname{Norm}_{x[i]}[\mu, \sigma^2] \right]$$

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



# 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.



Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince CS5340 :: Harold Soh

31

#### Example 1: Univariate Normal Distribution

#### Approach 1: Maximum Likelihood Estimation (MLE)

$$\hat{\mu}, \hat{\sigma}^2 = \underset{\mu, \sigma^2}{\operatorname{argmax}} \sum_{i=1}^{N} \log \left[ \operatorname{Norm}_{x[i]}[\mu, \sigma^2] \right]$$

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

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

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

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



#### Least Squares Interpretation

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

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

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

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

...gives `least squares' fitting criterion.

# Let's try it out.

#### https://github.com/crslab/CS5340-notebooks





CS5340 :: Harold Soh

#### MLE: Properties

- Easy and fast to compute
- Nice Asymptotic properties:
  - Consistent: if data generated from  $f(\theta^*)$ , MLE converges to its true value,  $\hat{\theta}_{MLE} \to \theta^*$  as  $n \to \infty$
  - 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  $\theta^*$ , and  $g(\theta^*)$  is a transformation of  $\theta^*$  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 <u>uncertainty</u> about the parameters estimates  $\mu$  and  $\sigma^2$ ?



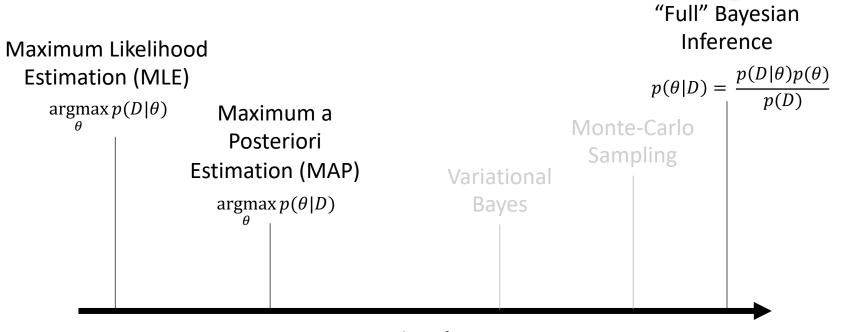


# Learning via Bayes

Bayesian Inference and Conjugate Models

## Learning Parameters

• Common approaches to learn the unknown parameters  $\theta$  from a set of given data  $\mathcal{D} = \{x[1], ..., x[N]\}$ :



#### **Computational Cost**

(In general and not to scale)



## 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



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

- Model uncertainty over  $\theta$  using prior distribution(s).
- Then, find posterior:

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

What can be a prior distribution?



Approach 3: Bayesian

Compute the posterior distribution using Bayes' rule:

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

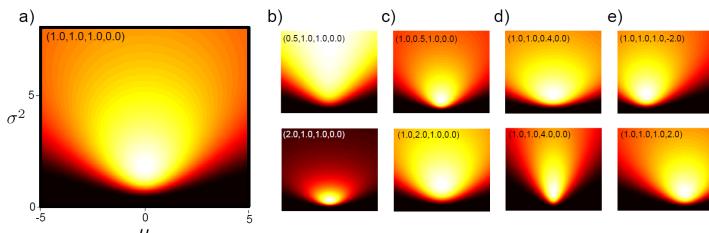
where:

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







CS5340 :: Harold Soh

### Normal Inverse Gamma

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



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  $p(\theta|D)$ ?

#### Approach 3: Bayesian

Compute the posterior distribution using Bayes' rule:

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

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

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

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)



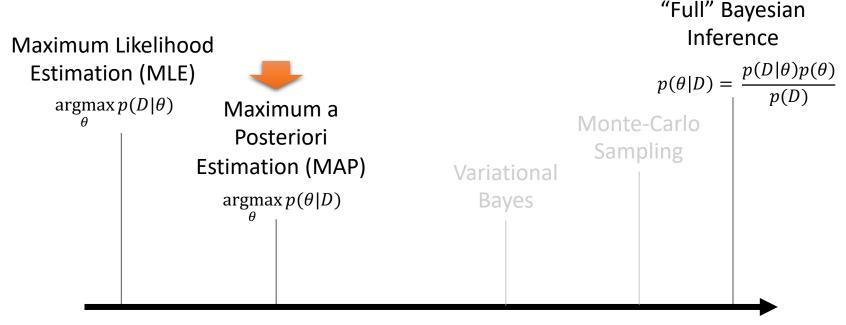


# Learning via MAP

Maximum a Posteriori Estimation(MAP)

## Learning Parameters

• Common approaches to learn the unknown parameters  $\theta$  from a set of given data  $\mathcal{D} = \{x[1], ..., x[N]\}$ :



**Computational Cost** 

(In general and not to scale)



## Maximum a Posteriori (MAP)

- Given data  $\mathcal{D} = \{x[1], ..., x[N]\}$
- Assume:
  - Joint distribution  $p(\mathcal{D}, \theta)$
  - Here  $\theta$  is a random variable
- Goal: Choose "good" heta
- The estimate  $\theta_{MAP}$  is a maximum aposteriori estimate (MAP) if

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



### Intuition: The "Peak" of the Posterior

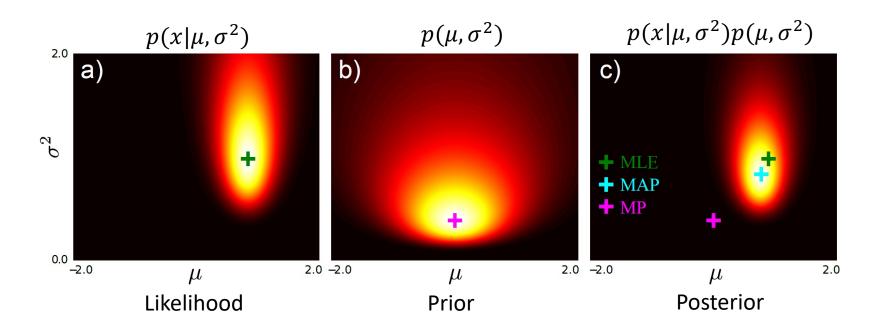




Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince CS5340:: Harold Soh 56

## Maximum a Posteriori (MAP)

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

$$\begin{split} &\theta_{MAP} = \underset{\theta}{\operatorname{argmax}}[p(\theta|D)] \\ &= \underset{\theta}{\operatorname{argmax}}\left[\frac{p(D|\theta)p(\theta)}{p(D)}\right] \qquad \text{(Bayes' rule)} \\ &= \underset{\theta}{\operatorname{argmax}}\left[\frac{\prod_{i=1}^{N}p(x[i]\mid\theta)\,p(\theta)}{p(D)}\right] \qquad \text{(i.i.d)} \\ &= \underset{\theta}{\operatorname{argmax}}\left[\prod_{i=1}^{N}p(x[i]\mid\theta)\,p(\theta)\right] \qquad \text{($p(D)$ is removed since it is independent of $\theta$)} \end{split}$$



Approach 2: Maximum a Posteriori (MAP)

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

Prior

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)



Approach 2: Maximum a Posteriori (MAP)

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

Prior

Likelihood: univariate Normal distribution

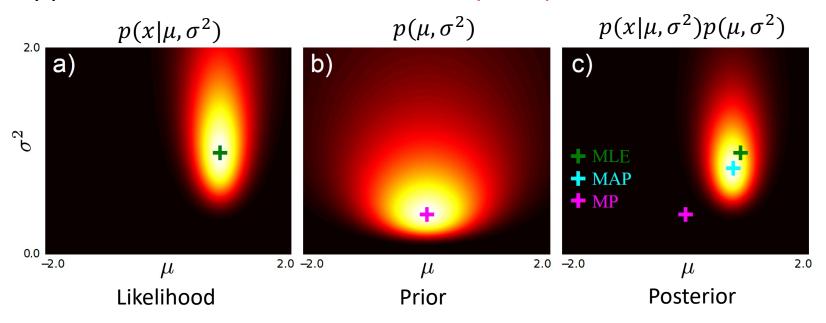
$$p(x|\mu,\sigma^2) = \prod_{i=1}^N \text{Norm}_{x[i]} [\mu,\sigma^2],$$

Prior: normal inverse gamma distribution

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



#### Approach 2: Maximum a Posteriori (MAP)



$$\hat{\mu}, \hat{\sigma}^{2} = \underset{\mu, \sigma^{2}}{\operatorname{argmax}} \left[ \prod_{i=1}^{N} p(x[i]|\mu, \sigma^{2}) p(\mu, \sigma^{2}) \right]$$

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



Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince CS5340 :: Harold Soh

60

#### Approach 2: Maximum a Posteriori (MAP)

$$\hat{\mu}, \hat{\sigma}^{2} = \underset{\mu, \sigma^{2}}{\operatorname{argmax}} \left[ \prod_{i=1}^{N} p(x[i] \mid \mu, \sigma^{2}) p(\mu, \sigma^{2}) \right]$$

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

#### Maximize the logarithm:

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



#### Approach 2: Maximum a Posteriori (MAP)

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

Taking derivatives and setting to zero:

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

We get:

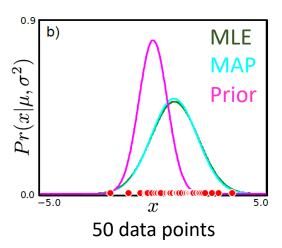
$$\hat{\mu} = \frac{\sum_{i} x[i] + \gamma \delta}{N + \gamma}, \qquad \hat{\sigma}^{2} = \frac{\sum_{i} (x[i] - \hat{\mu})^{2} + 2\beta + \gamma (\delta - \hat{\mu})^{2}}{N + 3 + 2\alpha}$$

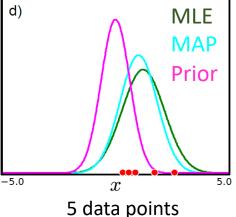
$$= \frac{N\bar{x} + \gamma \delta}{N + \gamma}$$



#### Approach 2: Maximum a Posteriori (MAP)

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





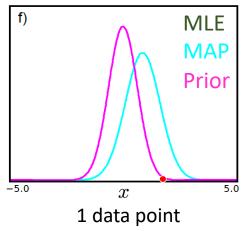




Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince

CS5340 :: Harold Soh 63

## Let's try it out



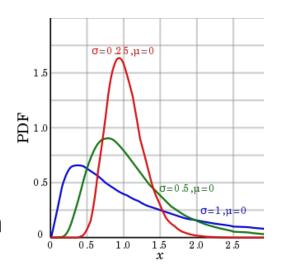


### Take-Home Exercise: Lognormal Prior

- Say you wanted to use a
  - normal prior for  $\mu$
  - lognormal prior for  $\sigma^2$
- Derive the MAP estimates
  - You can derive a closed-form solution for the mean
  - But would need to optimize for  $\sigma^2$



- Derive  $\mathcal{L} = \log p(D|\theta)p(\theta) = \log p(D|\theta) + \log p(\theta)$
- Then set  $\frac{\partial L}{\partial \theta_i} = 0$  for each parameter  $\theta_i$



## MAP: Properties

- Easy and fast to compute
- Incorporate prior information
- Avoid overfitting ("Regularization")
- As  $n \to \infty$ , 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  $\theta^*$ , and  $g(\theta^*)$  is a transformation of  $\theta^*$  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:**

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

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

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

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

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

## Bayesian Approach

#### **Predictive Density:**

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

Prediction for each possible heta

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

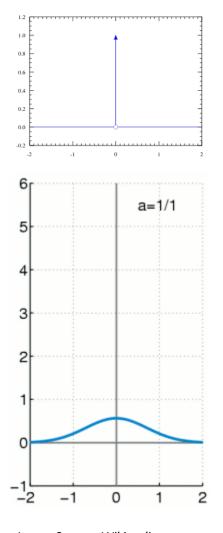


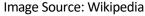
### Predictive Densities for 3 Approaches

How to rationalize different forms?

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

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







### Example 1: Univariate Normal Distribution

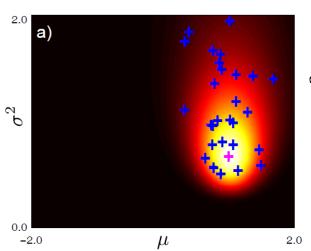
Approach 3: Bayesian

#### **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

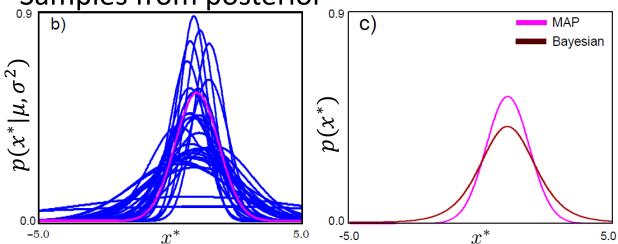




Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince CS5340:: Harold Soh 77

### Example 1: Univariate Normal Distribution

Approach 3: Bayesian

#### **Predictive density**

Take weighted sum of predictions from different parameter values:

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

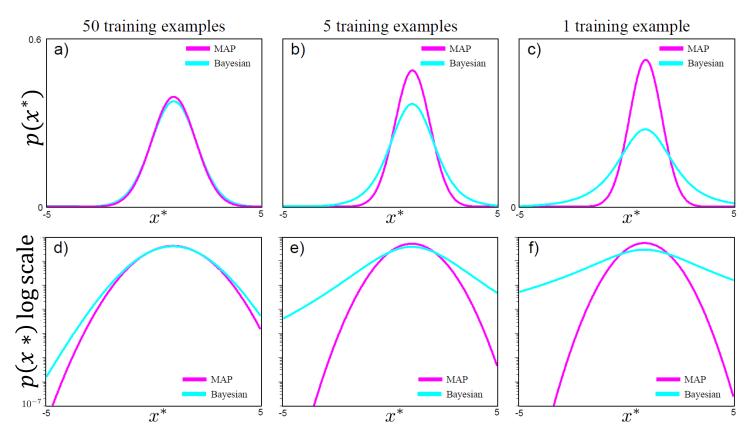
Where  $t_{2\widetilde{\alpha}}(x^*|\tilde{\delta},\frac{\widetilde{\beta}(\widetilde{\gamma}+1)}{\widetilde{\alpha}\widetilde{\gamma}})$  is the Generalized Student-T distribution with location  $\tilde{\delta}$  and scale  $\frac{\widetilde{\beta}(\widetilde{\gamma}+1)}{\widetilde{\alpha}\widetilde{\gamma}}$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



### 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.





 $Image\ Source: \ \hbox{``Computer Vision:}\ \ Models,\ Learning,\ and\ Inference'',\ Simon\ Prince$ 

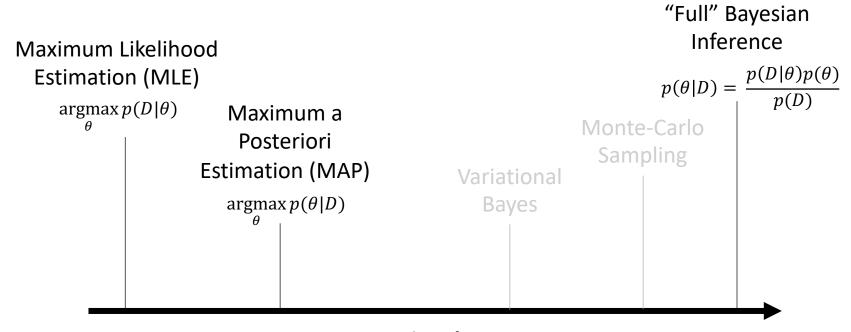
## Let's try it out





## Learning Parameters

• Common approaches to learn the unknown parameters  $\theta$  from a set of given data  $\mathcal{D} = \{x[1], ..., x[N]\}$ :



**Computational Cost** 

(In general and not to scale)





# Exponential Family

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

## **Exponential Family**

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

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

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



## Natural/Canonical form

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

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

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



## 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 $Z(\eta)$

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

Also called the normalizer:

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

Why? To get normalized distribution:

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

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

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

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

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



## The log Partition function $A(\eta)$

Alternatively, we can specify the log partition function:

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

• Is the log of the partition function:

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

• Why? To get normalized distribution for any  $\eta$ :

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



#### 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  $\eta_{MLE}$  satisfies:

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

- We can use this to solve for  $\eta_{MLE}$ 
  - 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!

• Poster But how can I find ExpFam distributions? distribution always has closed form solution (provided  $Z(\theta)$  is closed-form).



## Is Gaussian an ExpFam?

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

Rearrange to fit the ExpFam form:

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

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:  $A(\eta)$ 

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



## Many Distributions are ExpFam

**PMFs** 

**PDFs** 

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

- 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

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

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



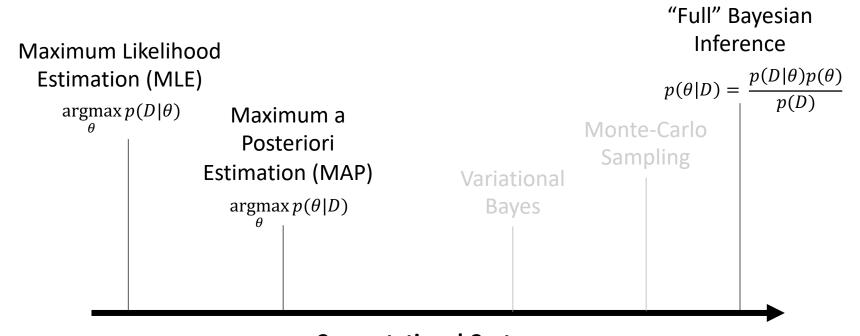


## Recap

MLE, MAP, Bayesian Inference, and Exponential Families

## Learning Parameters

• Common approaches to learn the unknown parameters  $\theta$  from a set of given data  $\mathcal{D} = \{x[1], ..., x[N]\}$ :



**Computational Cost** 

(In general and not to scale)



## **Exponential Family**

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

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

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



### Learning Outcomes

- Students should be able to:
  - 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.
  - 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

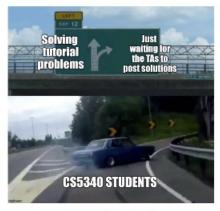


(c) Distracted Boyfriend





(b) Two Buttons Dilemma





#### CS5340 Meme of the Year

Model and learn parameters

$\overline{\mathbf{ID}}$	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

