

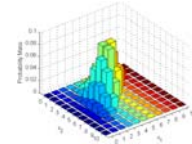
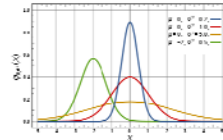
# CSE 446: Point Estimation Winter 2012

Dan Weld

Slides adapted from Carlos Guestrin (& Luke Zettlemoyer)

## Random Variable

- **Discrete or Continuous**
  - **Boolean:** like an atom in probabilistic logic
    - Denotes some event (eg, die will come up “6”)
  - **Continuous:** domain is R numbers
    - Denotes some quantity (age of a person taken from CSE446)
- **(Joint) Probability Distribution**



2

## Bayesian Methods

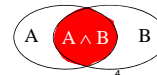
- Learning and classification methods based on probability theory.
  - Uses **prior** probability of each category  
Given no information about an item.
  - Produces a **posterior** probability distribution over possible categories  
Given a description of an item.
- Bayes theorem plays a critical role in probabilistic learning and classification.



3

## Axioms of Probability Theory

- All probabilities between 0 and 1  
 $0 \leq P(A) \leq 1$
- Probability of truth and falsity  
 $P(\text{true}) = 1 \quad P(\text{false}) = 0.$
- The probability of disjunction is:  
 $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$



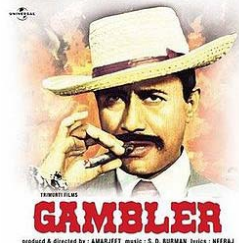
4

## De Finetti (1931)

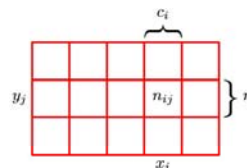
An agent who bets according to probabilities that violate these axioms...  
Can be forced to bet so as to **lose money**  
Regardless of the outcome



"Today's financial report will be a short one: We had money, now we don't."



## Terminology



**Marginal Probability**

$$p(X = x_i) = \frac{c_i}{N}$$

**Joint Probability**

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

**Conditional Probability**

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

↑  
X value is given

## Conditional Probability

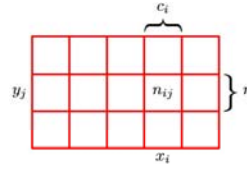
- $P(A | B)$  is the probability of  $A$  given  $B$
- Assumes:
  - $B$  is all and only information known.
- Defined by:

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$



7

## Probability Theory



Sum Rule

$$P(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij}$$

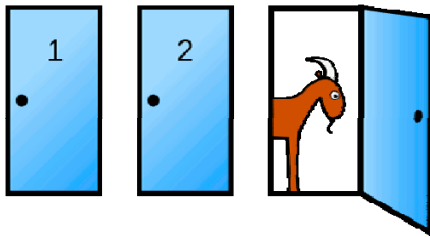
$$= \sum_{j=1}^L P(X = x_i, Y = y_j)$$

Product Rule

$$P(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N}$$

$$= P(Y = y_j | X = x_i) P(X = x_i)$$

## Monty Hall



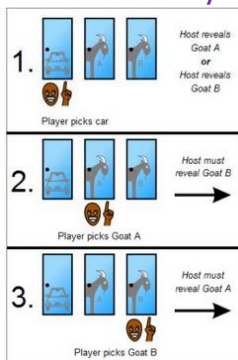
9

## Monty Hall



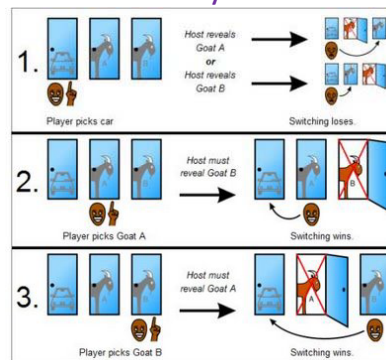
10

## Monty Hall



11

## Monty Hall



12

## Independence

- $A$  and  $B$  are *independent* iff:

$$P(A|B) = P(A) \quad \text{These constraints are logically equivalent}$$

$$P(B|A) = P(B)$$

- Therefore, if  $A$  and  $B$  are independent:

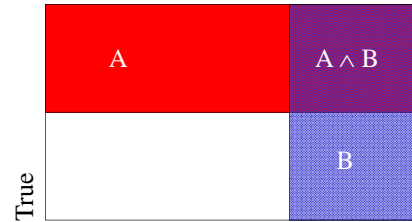
$$P(A|B) = \frac{P(A \wedge B)}{P(B)} = P(A)$$

$$P(A \wedge B) = P(A)P(B)$$

13

## Independence

$$P(A \wedge B) = P(A)P(B)$$

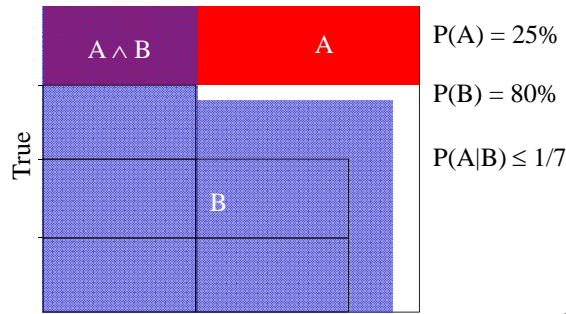


© Daniel S. Weld

14

## Independence Is Rare

$A$  &  $B$  *not* independent, since  $P(A|B) \neq P(A)$

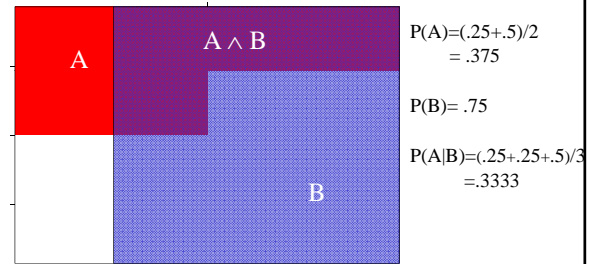


© Daniel S. Weld

15

## Conditional Independence....?

Are  $A$  &  $B$  independent?  $P(A|B) \leq P(A)$

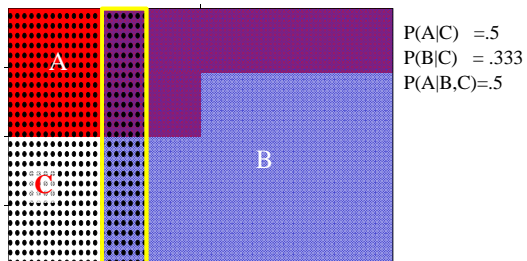


© Daniel S. Weld

16

## $A, B$ Conditionally Independent Given $C$

$$P(A|B, C) = P(A|C) \quad C = \text{spots}$$



© Daniel S. Weld

17

## Bayes Theorem



$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}$$

Simple proof from definition of conditional probability:

$$P(H|E) = \frac{P(H \wedge E)}{P(E)} \quad (\text{Def. cond. prob.})$$

$$\frac{P(H \wedge E)}{P(H)} = P(E|H) \quad (\text{Def. cond. prob.})$$

$$P(H \wedge E) = P(E|H)P(H) \quad (\text{Mult both sides of 2 by } P(H).)$$

**QED:**

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)} \quad (\text{Substitute 3 in 1.})$$

18

## Medical Diagnosis

*Encephalitis Computadoris*

Extremely Nasty

Extremely Rare 1/100M

DNA Testing 0% false negatives  
1/10,000 false positive



Scared?  $P(H|E) = \frac{P(E|H)P(H)}{P(E)}$

P(E)? In 100M people: 1+ 10,000

$P(H|E) = 1 * 0.00000001 / 0.00010001 \sim 0.000099$

## Your first consulting job

Billionaire (eccentric) Eastside tech founder asks:

- He says: I have thumbtack, if I flip it, what's the probability it will fall with the nail up?
- You say: Please flip it a few times:



- You say: The probability is:

- P(H) = 3/5

- He says: **Why???**

- You say: Because...

## Thumbtack – Binomial Distribution

- P(Heads) =  $\theta$ , P(Tails) =  $1-\theta$



- Flips are *i.i.d.*:
  - Independent events
  - Identically distributed according to Binomial distribution
- Sequence  $\mathcal{D}$  of  $\alpha_H$  Heads and  $\alpha_T$  Tails

$$P(\mathcal{D} | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

## Maximum Likelihood Estimation

- **Data:** Observed set  $\mathcal{D}$  of  $\alpha_H$  Heads and  $\alpha_T$  Tails
- **Hypothesis:** Binomial distribution
- **Learning:** finding  $\theta$  is an optimization problem
  - What's the objective function?

$$P(\mathcal{D} | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

- **MLE:** Choose  $\theta$  to maximize probability of  $\mathcal{D}$

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} P(\mathcal{D} | \theta) \\ &= \arg \max_{\theta} \ln P(\mathcal{D} | \theta) \end{aligned}$$

## Your first parameter learning algorithm

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} \ln P(\mathcal{D} | \theta) \\ &= \arg \max_{\theta} \ln \theta^{\alpha_H} (1 - \theta)^{\alpha_T} \end{aligned}$$

- Set derivative to zero, and solve!

$$\begin{aligned} \frac{d}{d\theta} \ln P(\mathcal{D} | \theta) &= \frac{d}{d\theta} [\ln \theta^{\alpha_H} (1 - \theta)^{\alpha_T}] \\ &= \frac{d}{d\theta} [\alpha_H \ln \theta + \alpha_T \ln(1 - \theta)] \\ &= \alpha_H \frac{d}{d\theta} \ln \theta + \alpha_T \frac{d}{d\theta} \ln(1 - \theta) \\ &= \frac{\alpha_H}{\theta} - \frac{\alpha_T}{1 - \theta} = 0 \end{aligned}$$

$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

## But, how many flips do I need?

$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

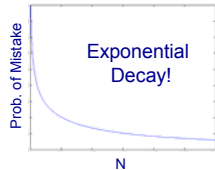
- Billionaire says: I flipped 3 heads and 2 tails.
- You say:  $\theta = 3/5$ , I can prove it!
- He says: What if I flipped 30 heads and 20 tails?
- You say: Same answer, I can prove it!
- He says: **What's better?**
- You say: Umm... The more the merrier???
- He says: Is this why I am paying you the big bucks???

## A bound (from Hoeffding's inequality)

For  $N = \alpha_H + \alpha_T$  and  $\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$

Let  $\theta^*$  be the true parameter, for any  $\epsilon > 0$ :

$$P(|\hat{\theta} - \theta^*| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$$



## PAC Learning

- PAC: Probably Approximate Correct
- Billionaire says: I want to know the thumbtack  $\theta$ , within  $\epsilon = 0.1$ , with probability at least  $1 - \delta = 0.95$ .
- How many flips? Or, how big do I set  $N$ ?

$$P(|\hat{\theta} - \theta^*| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$$

$$\delta \geq 2e^{-2N\epsilon^2} \geq P(\text{mistake})$$

$$\ln \delta \geq \ln 2 - 2N\epsilon^2$$

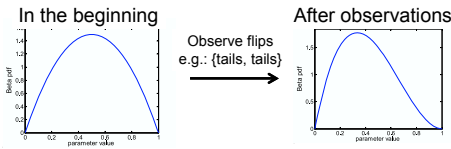
$$N \geq \frac{\ln(2/\delta)}{2\epsilon^2}$$

Interesting! Lets look at some numbers!  
 $\epsilon = 0.1, \delta = 0.05$

$$N \geq \frac{\ln(2/0.05)}{2 \times 0.1^2} \approx \frac{3.8}{0.02} = 190$$

## What if I have prior beliefs?

- Billionaire says: Wait, I know that the thumbtack is "close" to 50-50. What can you do for me now?
- You say: I can learn it the Bayesian way...
- Rather than estimating a single  $\theta$ , we obtain a distribution over possible values of  $\theta$



## Bayesian Learning

Use Bayes rule:

$$P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta)P(\theta)}{P(\mathcal{D})}$$

Labels: Posterior (points to the result), Data Likelihood (points to the numerator), Prior (points to the prior distribution graph), Normalization (points to the denominator).

Or equivalently:  $P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta)$

Also, for uniform priors:

→ reduces to MLE objective

$$P(\theta) \propto 1 \quad P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)$$

## Bayesian Learning for Thumbtacks

$$P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta)$$

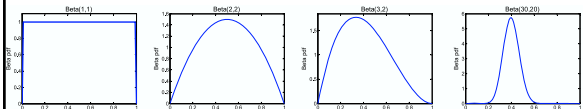
Likelihood function is Binomial:

$$P(\mathcal{D} | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

- What about prior?
  - Represent expert knowledge
  - Simple posterior form
- Conjugate priors:
  - Closed-form representation of posterior
  - For Binomial, conjugate prior is Beta distribution

## Beta prior distribution – P(θ)

$$P(\theta) = \frac{\theta^{\beta_H - 1} (1 - \theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim \text{Beta}(\beta_H, \beta_T)$$



• Likelihood function:  $P(\mathcal{D} | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$

• Posterior:  $P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta)$

$$P(\theta | \mathcal{D}) \propto \theta^{\alpha_H} (1 - \theta)^{\alpha_T} \theta^{\beta_H - 1} (1 - \theta)^{\beta_T - 1}$$

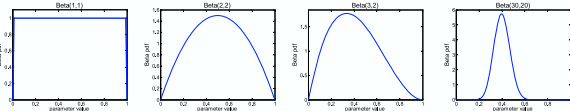
$$= \theta^{\alpha_H + \beta_H - 1} (1 - \theta)^{\alpha_T + \beta_T - 1}$$

$$= \text{Beta}(\alpha_H + \beta_H, \alpha_T + \beta_T)$$

## Posterior Distribution

- Prior:  $Beta(\beta_H, \beta_T)$
- Data:  $\alpha_H$  heads and  $\alpha_T$  tails
- Posterior distribution:  

$$P(\theta | \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$



## Bayesian Posterior Inference

- Posterior distribution:

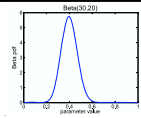
$$P(\theta | \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

- Bayesian inference:

- No longer single parameter
- For any specific  $f$ , the function of interest
- Compute the expected value of  $f$

$$E[f(\theta)] = \int_0^1 f(\theta) P(\theta | \mathcal{D}) d\theta$$

- Integral is often hard to compute



## MAP: Maximum a Posteriori Approximation

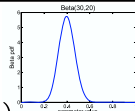
$$P(\theta | \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

$$E[f(\theta)] = \int_0^1 f(\theta) P(\theta | \mathcal{D}) d\theta$$

- As more data is observed, Beta is more certain
- MAP: use most likely parameter to approximate the expectation

$$\hat{\theta} = \arg \max_{\theta} P(\theta | \mathcal{D})$$

$$E[f(\theta)] \approx f(\hat{\theta})$$



## MAP for Beta distribution

$$P(\theta | \mathcal{D}) = \frac{\theta^{\beta_H + \alpha_H - 1} (1 - \theta)^{\beta_T + \alpha_T - 1}}{B(\beta_H + \alpha_H, \beta_T + \alpha_T)} \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

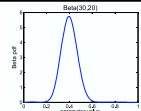
- MAP: use most likely parameter:

$$\hat{\theta} = \arg \max_{\theta} P(\theta | \mathcal{D}) = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

Beta prior equivalent to extra thumbtack flips

As  $N \rightarrow \infty$ , prior is "forgotten"

**But, for small sample size, prior is important!**



## Envelope Problem

One envelope has twice as much money as other

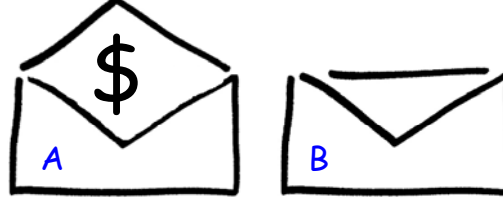


Switch?

35

## Envelope Problem

One envelope has twice as much money as other



A has \$20

$$E[B] = \$ \frac{1}{2} (10 + 40)$$

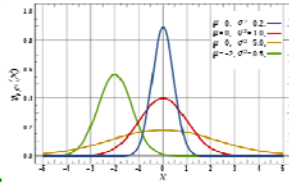
$$= \$ 25$$

Switch?

36

## What about continuous variables?

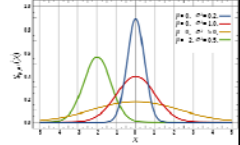
- Billionaire says: If I am measuring a continuous variable, what can you do for me?
- You say: Let me tell you about Gaussians...



$$P(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

## Some properties of Gaussians

- Affine transformation (multiplying by scalar and adding a constant) are Gaussian
  - $X \sim N(\mu, \sigma^2)$
  - $Y = aX + b \rightarrow Y \sim N(a\mu + b, a^2\sigma^2)$
- Sum of Gaussians is Gaussian
  - $X \sim N(\mu_x, \sigma_x^2)$
  - $Y \sim N(\mu_y, \sigma_y^2)$
  - $Z = X + Y \rightarrow Z \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$



- Easy to differentiate, as we will see soon!

## Learning a Gaussian

- Collect a bunch of data
  - Hopefully, i.i.d. samples
  - e.g., exam scores
- Learn parameters
  - Mean:  $\mu$
  - Variance:  $\sigma$

$X_i$ $i =$	Exam Score
0	85
1	95
2	100
3	12
...	...
99	89

$$P(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

## MLE for Gaussian: $P(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

- Prob. of i.i.d. samples  $D = \{x_1, \dots, x_N\}$ :

$$P(\mathcal{D} | \mu, \sigma) = \left( \frac{1}{\sigma\sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\mu_{MLE}, \sigma_{MLE} = \arg \max_{\mu, \sigma} P(\mathcal{D} | \mu, \sigma)$$

- Log-likelihood of data:

$$\begin{aligned} \ln P(\mathcal{D} | \mu, \sigma) &= \ln \left[ \left( \frac{1}{\sigma\sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right] \\ &= -N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

## Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for mean?

$$\begin{aligned} \frac{d}{d\mu} \ln P(\mathcal{D} | \mu, \sigma) &= \frac{d}{d\mu} \left[ N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\ &= \frac{d}{d\mu} \left[ -N \ln \sigma\sqrt{2\pi} \right] - \sum_{i=1}^N \frac{d}{d\mu} \left[ \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\ &= - \sum_{i=1}^N \frac{(x_i - \mu)}{\sigma^2} = 0 \\ &= - \sum_{i=1}^N x_i + N\mu = 0 \end{aligned}$$

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

## MLE for variance

- Again, set derivative to zero:

$$\begin{aligned} \frac{d}{d\sigma} \ln P(\mathcal{D} | \mu, \sigma) &= \frac{d}{d\sigma} \left[ -N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\ &= \frac{d}{d\sigma} \left[ -N \ln \sigma\sqrt{2\pi} \right] - \sum_{i=1}^N \frac{d}{d\sigma} \left[ \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\ &= -\frac{N}{\sigma} + \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^3} = 0 \end{aligned}$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

## Learning Gaussian parameters

- MLE:

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

- BTW. MLE for the variance of a Gaussian is **biased**

– Expected result of estimation is **not** true parameter!

– Unbiased variance estimator:

$$\hat{\sigma}_{unbiased}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

## Bayesian learning of Gaussian parameters

- Conjugate priors

– Mean: Gaussian prior

– Variance: Wishart Distribution

- Prior for mean:

$$P(\mu | \eta, \lambda) = \frac{1}{\lambda\sqrt{2\pi}} e^{-\frac{(\mu-\eta)^2}{2\lambda^2}}$$

## MAP for mean of Gaussian

$$P(\mu | \eta, \lambda) = \frac{1}{\lambda\sqrt{2\pi}} e^{-\frac{(\mu-\eta)^2}{2\lambda^2}} \quad P(\mathcal{D} | \mu, \sigma) = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^N \prod_{i=1}^N e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}$$

$$\frac{d}{d\mu} [\ln P(\mathcal{D} | \mu) P(\mu)] = \frac{d}{d\mu} [\ln P(\mathcal{D} | \mu) + \ln P(\mu)]$$

$$\begin{aligned} \frac{\partial}{\partial \mu} \ln P(\mathcal{D} | \mu) &= \sum_i \frac{(x_i - \mu)}{\sigma^2} \\ \frac{\partial}{\partial \mu} \ln \left[ \frac{1}{\lambda\sqrt{2\pi}} e^{-\frac{(\mu-\eta)^2}{2\lambda^2}} \right] &= -\frac{(\mu-\eta)}{\lambda^2} \end{aligned} \quad \left| \quad \begin{aligned} \sum_i \frac{(x_i - \mu)}{\sigma^2} + \frac{(\eta - \mu)}{\lambda^2} &= 0 \\ \frac{N\mu}{\sigma^2} + \frac{\mu}{\lambda^2} &= \frac{\sum_i x_i}{\sigma^2} + \frac{\eta}{\lambda^2} \\ \Rightarrow \mu &= \frac{1}{\frac{N}{\sigma^2} + \frac{1}{\lambda^2}} \left[ \frac{\sum_i x_i}{\sigma^2} + \frac{\eta}{\lambda^2} \right] \end{aligned} \right.$$