

CS5340: Tutorial 2

Asst. Prof. Harold Soh
TA: Eugene Lim



Poll Everywhere!

<https://pollev.com/haroldsohsoo986>



Upcoming Dates

- **Project**

- Form teams: Feb 6th
- Abstract due: Feb 27th

- Piazza teams

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 + Langevin Dynamics
11	2 Apr	Inference and Decision-Making	Diffusion Models + Sequential VAEs
12	9 Apr	Gaussian Processes (optional)	Quiz 2
13	16 Apr	Project Presentations	Closing Lecture

CS5340: Tutorial 2

Asst. Prof. Harold Soh

TA: Eugene Lim

1. Uncorrelated Random Variables



a. Two random vars X and Y where

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] = 0$$

<https://pollev.com/haroldsohsoo986>

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

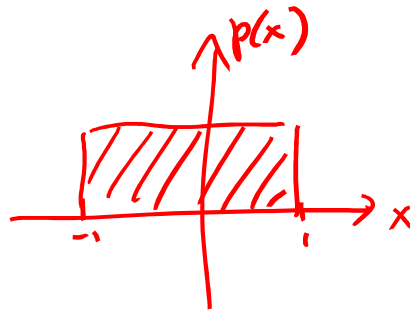
Are X and Y **independent**? Justify your answer.

if X and Y are independent $\Rightarrow \text{Cov}[X, Y] = 0$



$$X \sim \text{Uniform}[-1, 1]$$

$$Y = X^2$$

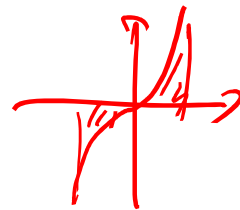
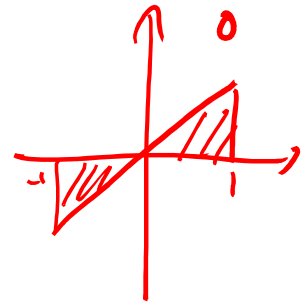


$$\text{Cov}[X, Y] = E[XY] - E[X]E[Y].$$

$$\Rightarrow E[\underline{X^3}] - E[\underline{X}]E[X^2] = 0.$$

$$E[X] = \int_{-1}^1 x p(x) dx = \int_{-1}^1 x \frac{1}{2} dx = 0$$

$$E[X^3] = \int_{-1}^1 x^3 p(x) dx = \int_{-1}^1 \frac{x^3}{2} dx = 0$$



1. Uncorrelated Random Variables

b. Consider:

- samples x_1, \dots, x_N drawn from $p(X)$
- samples y_1, \dots, y_N drawn from $p(Y)$

We want to model the joint $p_\theta(X, Y)$

Can we just perform MLE by finding

$$\operatorname{argmax}_\theta \sum_i^N \log p_\theta(x_i, y_i) ?$$



<https://pollev.com/haroldsohsoo986>

assume $Y = X$

X	Y
0	0
1	1
0	0
1	1
1	1
0	0

$X \sim \text{Bern}(0.5)$

$p(x, y)$

$(x, y) \sim p(x, y)$

x	y	
0	0	$\theta_{00} = 0.5$
0	1	$\theta_{01} = 0$
1	0	$\theta_{10} = 0$
1	1	$\theta_{11} = 0.5$

$X \sim p(x)$

$Y \sim p(y)$

x	y
0	0
1	0
0	0
0	0
1	0
1	0
0	0
0	0
1	0

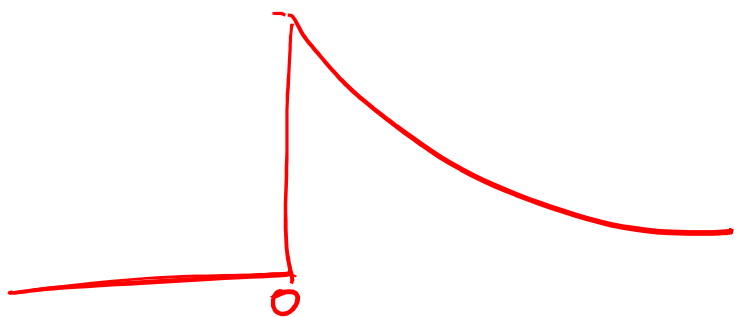
x	y	
0	0	0.25
0	1	0.25
1	0	0.25
1	1	0.25

Exponential Family (ExpFam)

2.a. Show the exponential distribution is exponential family (ExpFam).

$$p(x) = \frac{h(x) \exp(\eta^T s(x))}{Z(\eta)} = h(x) \exp(\eta^T s(x) - A(\eta)) \quad A(\eta) = \log Z(\eta)$$

$$p(x) = 1 \times \underline{\lambda \exp(-\lambda x)} \quad \underline{x \geq 0}.$$



$$\begin{aligned} s(x) &= x \\ \eta &= -\lambda \\ \rightarrow h(x) &= 1 \text{ if } x \geq 0, \quad 0 \text{ otherwise} \\ Z(\eta) &= \frac{1}{\lambda} = \frac{-1}{\eta} \\ A(\eta) &= \log\left(\frac{-1}{\eta}\right) \end{aligned}$$

ExpFam MLE

$$\mathbb{E}_{x \sim D(\eta)} [s(x)] = \nabla_{\eta} A(\eta) \quad \text{--- ①}$$

$$\nabla_{\eta} A(\eta_{MLE}) = \frac{1}{N} \sum_{i=1}^N x_i, \quad x_1, \dots, x_N \quad \text{--- ②}$$

2.b. Derive the MLE for the exponential distribution using facts about ExpFam.

1. Find $\nabla_{\eta} A(\eta)$

2. Evaluate $\nabla_{\eta} A(\eta)$ at η_{MLE} .

3. Use ② and solve.

1. Find $\nabla_{\eta} A(\eta)$:

$$\frac{d}{d\eta} \log\left(\frac{1}{\eta}\right) = \frac{1/\eta^2}{(-1/\eta)} = -\frac{1}{\eta}$$

2. Evaluate $\dots \eta_{MLE}$, 3. Use ② to solve

$$\frac{-1}{\eta_{MLE}} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\Rightarrow \eta_{MLE} = \frac{-1}{\frac{1}{N} \sum_{i=1}^N x_i}$$

$$\lambda_{MLE} = -\eta_{MLE} = \frac{1}{\frac{1}{N} \sum_{i=1}^N x_i}$$

Gaussian is ExpFam

$$p(x) = \underline{h(x)} \exp(\underline{\eta^T s(x)} - \underline{A(\eta)})$$

$$h(x) = 1 \quad A(\eta) = \sqrt{2\pi}\sigma^2$$

2.c.1. Show that the Gaussian is ExpFam.

$$\begin{aligned} p(x) &= \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(\frac{-(x-\mu)^2}{2\sigma^2}\right) \leftarrow \\ &= \frac{1}{\sqrt{2\pi}} \exp(\log \frac{1}{\sigma}) \exp\left(\frac{-(x^2 - 2x\mu + \mu^2)}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(\log \frac{1}{\sigma} - \frac{x^2}{2\sigma^2} + \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(\underbrace{\frac{\mu}{\sigma^2} x - \frac{1}{2\sigma^2} x^2}_{\eta^T s(x)} - \left(\log \sigma + \frac{\mu^2}{2\sigma^2}\right)\right) \\ &= \end{aligned}$$

$$\begin{aligned} &\left[\begin{array}{c} \mu/\sigma^2 \\ -1/2\sigma^2 \end{array} \right]^T \left[\begin{array}{c} x \\ x^2 \end{array} \right] \\ &\eta \quad s(x) \end{aligned}$$

$$\left\{ \begin{aligned} h(x) &= \frac{1}{\sqrt{2\pi}} \\ A(\eta) &= \log \sigma + \frac{\mu^2}{2\sigma^2} \\ &= \frac{-\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2) \\ \eta &= \left[\begin{array}{c} \mu/\sigma^2 \\ -1/2\sigma^2 \end{array} \right] \leftarrow \begin{array}{l} \eta_1 \\ \eta_2 \end{array} \\ s(x) &= \left[\begin{array}{c} x \\ x^2 \end{array} \right] \end{aligned} \right.$$

try to
verify it
yourself

$$\hat{\eta}_1 : \text{MLE}$$

Gaussian MLE

$$\nabla_{\eta} A(\eta_{\text{MLE}}) = \frac{1}{N} \sum_{i=1}^N s(x_i)$$

$$s(x_i) = \begin{bmatrix} x_i \\ x_i^2 \end{bmatrix}$$

$$A(\eta) = \frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2)$$

2.c.2. Derive the MLE of the Gaussian's natural parameters.

$$\nabla_{\eta} A(\eta) = \begin{bmatrix} \frac{\partial}{\partial \eta_1} A(\eta) \\ \frac{\partial}{\partial \eta_2} A(\eta) \end{bmatrix}$$

$$\frac{\partial}{\partial \eta_1} A(\eta) = \frac{-2\eta_1}{4\eta_2} = \frac{-\eta_1}{2\eta_2}$$

$$\begin{aligned} \frac{\partial}{\partial \eta_2} A(\eta) &= \frac{\eta_1^2}{4\eta_2^2} - \frac{1}{2} \left(\frac{-2}{-2\eta_2} \right) \\ &= \frac{\eta_1^2}{4\eta_2^2} - \frac{1}{2\eta_2} \end{aligned}$$

$$\begin{aligned} \frac{-2\hat{\eta}_1}{4\hat{\eta}_2} &= \frac{1}{N} \sum_{i=1}^N x_i \\ &= \hat{\mu} \end{aligned}$$

$$\frac{\frac{\hat{\eta}^2}{\hat{\lambda}^2}}{4\hat{\eta}^2} - \frac{1}{2\hat{\eta}^2} = \frac{1}{N} \sum_{i=1}^N x_i^2$$

$= \hat{\mu}^2 + \hat{\sigma}^2$

$$\frac{\frac{\hat{\eta}^2}{\hat{\lambda}^2}}{4\hat{\eta}^2} = \frac{\left(\frac{\hat{\mu}}{\hat{\sigma}^2}\right)^2}{4\left(\frac{-1}{2\hat{\sigma}^2}\right)^2} = \hat{\mu}^2$$

$$+ \frac{-1}{2\hat{\eta}^2} = \frac{-1}{2\left(\frac{-1}{2\hat{\sigma}^2}\right)} = \hat{\sigma}^2$$

$$\hat{\mu}^2 + \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N x_i^2$$

$$\hat{\sigma}^2 = \left(\frac{1}{N} \sum_{i=1}^N x_i^2 \right) - \frac{\hat{\mu}^2 \cdot N}{N}$$

$\sum \frac{1}{N} \hat{\mu}^2$

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

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

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

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

$$\frac{1}{N} \sum x_i \hat{\mu}$$

$$= \frac{1}{N} \hat{\mu} \sum x_i$$

$$= \hat{\mu} \hat{\mu}$$

$$= \hat{\mu}^2$$

Questions?

<https://pollev.com/haroldsohsoo986>



3. 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



<https://pollev.com/haroldsohsoo986>

3. Meme of the Year

a. What model (distribution) can we use to model this data? Compute the parameters via MLE.

$$x_1 = [1 \ 0 \ 0 \ 0]$$

$$x_2 = [0 \ 0 \ 1 \ 0]$$

Categorical Distribution

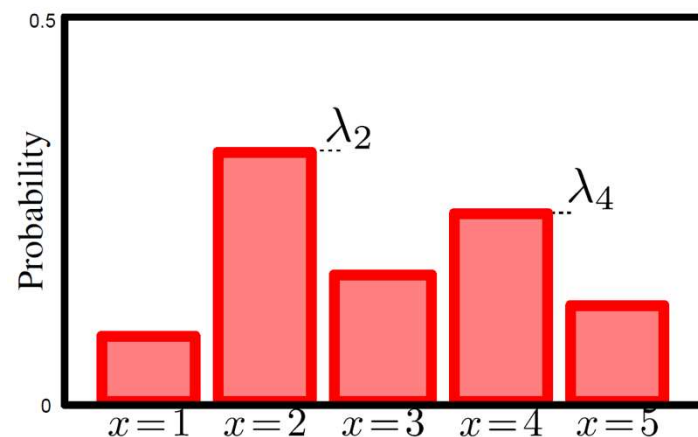
- Discrete variables X that take on **1-of- K possible mutually exclusive states**, e.g. a K -faced die.
- \mathbf{x} is represented by a **K -dimensional vector** \mathbf{e}_k in which one of the elements $x_k = 1$, and $\sum_{k=1}^K x_k = 1$.
- e.g. $K = 5$, and $\mathbf{x} = \mathbf{e}_3 = [0, 0, 1, 0, 0]^T$.
- **K parameters** $\lambda = [\lambda_1, \dots, \lambda_K]^T$, where $\lambda \geq 0$, $\sum_k \lambda_k = 1$.

$$p(X = \mathbf{e}_k | \lambda) = \lambda_k$$

Or

$$p(\mathbf{x}) = \prod_{k=1}^K \lambda_k^{x_k} = \lambda_k,$$

$$p(\mathbf{x}) = \text{Cat}_x[\lambda] \quad ||$$



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

3. Meme of the Year

a. What model (distribution) can we use to model this data? Compute ~~ε~~ the parameters via MLE.

$$\begin{aligned}
 & \text{Cat}[\underline{\vec{\lambda}}] \quad D = \{x_1, \dots, x_N\}. \\
 & \boxed{\vec{\lambda}_{MLE} = \underset{\vec{\lambda}}{\operatorname{argmax}} p(D|\vec{\lambda})} \\
 & p(D|\vec{\lambda}) = \prod_{i=1}^N \prod_{j=1}^k \lambda_j^{[x_{ij}=1]} \quad \left(\sum_{j=1}^k \lambda_j = 1 \right) \\
 & \quad \quad \quad = \prod_{j=1}^k \lambda_j^{N_j} \\
 & \underset{\vec{\lambda}}{\operatorname{argmax}} \log p(D|\vec{\lambda}) = \sum_{j=1}^k N_j \log \lambda_j + v \left(\sum_{j=1}^k \lambda_j - 1 \right).
 \end{aligned}$$

maximize $f(\theta)$
s.t. $g(\theta) = 0$

$$\mathcal{L} = \sum_{j=1}^k N_j \log \lambda_j + v \left(\sum_{j=1}^k \lambda_j - 1 \right)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_j} = \frac{N_j}{\lambda_j} + v = 0 \Rightarrow \lambda_j = -\frac{N_j}{v}$$

Const.: $\sum_{j=1}^k \lambda_j = 1 \Rightarrow \sum_{j=1}^k -\frac{N_j}{v} = 1 \Rightarrow v = -\sum_{j=1}^k N_j$

$$\lambda_j = \frac{N_j}{\sum_{l=1}^k N_l}$$

3. Meme of the Year

b. What prior distribution can we use for our model?

$$p(\vec{\lambda})$$

Dirichlet Distribution

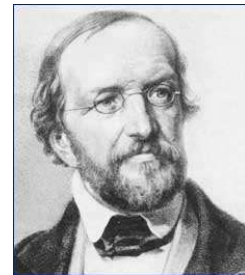
$$p(x) = \frac{1}{Z} \prod_{i=1}^K \lambda_i^{x_i-1}$$

$\sum \lambda_i = 1$
 $[\lambda_1, \lambda_2, \lambda_3]$

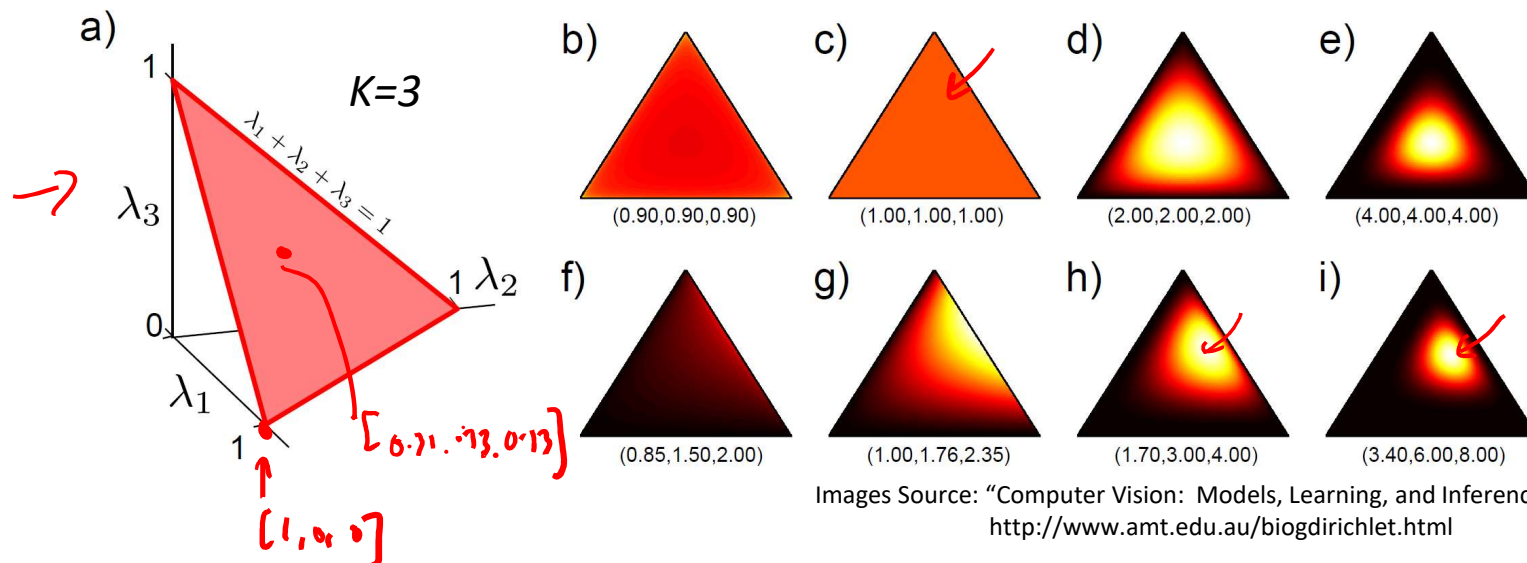
- Conjugate distribution of **categorical distribution**.
- Defined over K parameters of Categorical distribution, $\lambda_k \in [0,1]$, where $\sum_k \lambda_k = 1$.

$$p(\lambda_1, \dots, \lambda_K) = \frac{\Gamma[\sum_{k=1}^K \alpha_k]}{\prod_{k=1}^K \Gamma[\alpha_k]} \prod_{k=1}^K \lambda_k^{\alpha_k-1},$$

$$p(\lambda_1, \dots, \lambda_K) = \text{Dir}_{\lambda_1 \dots \lambda_K}[\alpha_1, \dots, \alpha_K]$$



Peter Gustav Lejeune Dirichlet
(1805-1859)



Dirichlet Distribution

Useful property:
 $\Gamma(1) = 1$
 $\Gamma(z + 1) = z\Gamma(z)$

$$p(\lambda_1, \dots, \lambda_K) = \frac{\Gamma[\sum_{k=1}^K \alpha_k]}{\prod_{k=1}^K \Gamma[\alpha_k]} \prod_{k=1}^K \lambda_k^{\alpha_k - 1},$$

$$p(\lambda_1, \dots, \lambda_K) = \text{Dir}_{\lambda_1 \dots \lambda_K}[\alpha_1, \dots, \alpha_K]$$

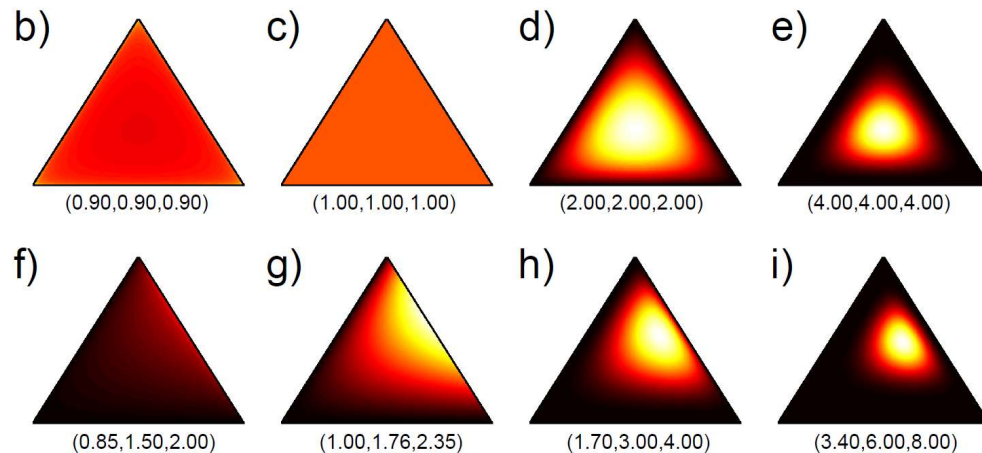
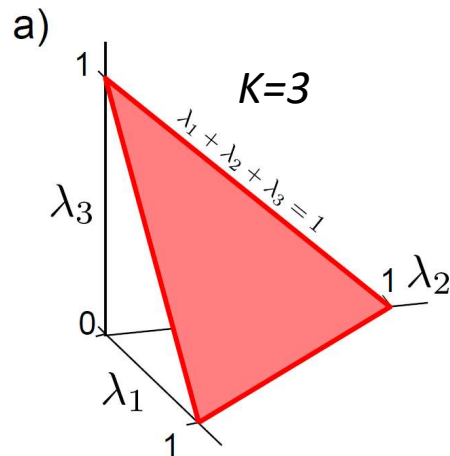
Gamma & Beta Functions:

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$$

$\rightarrow \Gamma(n) = (n-1)!$

$$\rightarrow B(\underline{\alpha}) = \frac{\prod_{k=1}^K \Gamma[\alpha_k]}{\Gamma[\sum_{k=1}^K \alpha_k]}$$

- K hyperparameters $\alpha_k > 0$.



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

3. Meme of the Year

b. What is the posterior distribution?

$$p(\vec{\lambda}) = \frac{1}{B(\vec{\alpha})} \prod_{j=1}^k \lambda_j^{\alpha_j - 1}$$

$$p(x | \vec{\lambda}) = \prod_{j=1}^k \lambda_j^{N_j}$$

Bayes

$$p(\vec{\lambda} | x) \propto p(x | \vec{\lambda}) p(\vec{\lambda})$$

$$\left[\prod_{j=1}^k \lambda_j^{N_j} \right] \left[\frac{1}{B(\vec{\alpha})} \prod_{j=1}^k \lambda_j^{\alpha_j - 1} \right]$$

$$= \underbrace{\frac{1}{B(\vec{\alpha})}}_{\text{const w.r.t } \lambda} \prod_{j=1}^k \lambda_j^{\alpha_j + N_j - 1} \Rightarrow \text{Dir} \left[\begin{matrix} \vec{\alpha} \\ n \end{matrix} \right]$$

$$[\alpha_1 + N_1, \alpha_2 + N_2, \dots, \alpha_k + N_k]^T$$

3. Meme of the Year

Useful property:

$$\Gamma(1) = 1$$

$$\Gamma(z + 1) = z\Gamma(z)$$

b. What is the posterior predictive distribution?

$$p(\hat{\theta} | X)$$

$$p(x^* | X) = \int \overbrace{p(x^* | \vec{\lambda})}^{\text{Co + likelihood}} \overbrace{p(\vec{\lambda} | X)}^{\text{Dirich}[\vec{\alpha}]}} d\lambda$$

$$= \int \prod_{j=1}^k \lambda_j^{[x_j^*=1]} \cdot \frac{1}{B(\vec{\alpha})} \prod_{j=1}^k \lambda_j^{\tilde{\alpha}_j-1} d\lambda.$$

multiply and divide by

$$\frac{\prod_{j=1}^k \Gamma(\tilde{\alpha}_j + [x_j^*=1])}{\Gamma(\sum_{j=1}^k \tilde{\alpha}_j + [x_j^*=1])}$$

Questions?

<https://pollev.com/haroldsohsoo986>



Please Prepare for Next Week

- Watch the videos
 - Bayesian Networks!
- Do the tutorial

