3.2) With rule of chain in probability

$$P(x_{1:N}) = P(x_1) P(x_2/x_1) P(x_3/x_{1:2}) \cdots$$

After seeing dataset $D = H, T, T, H, H$ and if we
have prior probabilities as $P(T) = \alpha_0$, $P(H) = \alpha_1$
and $\alpha = \alpha_0 + \alpha_1$. Then we can update derive
likelihood by

$$P(D) = P(H) P(T/H) P(T/_{H,T}) P(^H/_{H,T,T}) P(^H/_{H,T,T,H})$$

$$= \frac{\alpha_1}{\alpha} \cdot \frac{\alpha_0}{\alpha+1} \cdot \frac{\alpha_0+1}{\alpha+2} \cdot \frac{\alpha_1+1}{\alpha+3} \cdot \frac{\alpha_1+2}{\alpha+4}.$$

$$= \frac{[(\alpha_0)(\alpha_0+1)][(\alpha_1)(\alpha_1+1)(\alpha_1+2)]}{[\alpha \cdot (\alpha+1) \cdot (\alpha+2)(\alpha+3)(\alpha+4)]}$$

To generalize

$$= \frac{[(\alpha_0)(\alpha_0+1) \cdots (\alpha_0+N_0-1)] [(\alpha_1)(\alpha_1+1) \cdots (\alpha_1+N_1-1)]}{[(\alpha)(\alpha+1) \cdots (\alpha+N-1)]}$$

We can rewrite this as

$$(\because \Gamma(\alpha) = (\alpha-1)!)$$

$$= \frac{\Gamma(\alpha_0+N_0)}{\Gamma(\alpha_0)} \cdot \frac{\Gamma(\alpha_1+N_1)}{\Gamma(\alpha_1)} \bigg/ \frac{\Gamma(\alpha+N)}{\Gamma(\alpha)}$$

$$= \frac{\Gamma(\alpha_0+N_0) \Gamma(\alpha_1+N_1)}{\Gamma(\alpha_1+\alpha_0+N)} \cdot \frac{\Gamma(\alpha_1+\alpha_0)}{\Gamma(\alpha_0)\Gamma(\alpha_1)} \bigg/\!\!/$$

Beta-binomial posterior predictive

$$P(x/y,D) = \int Bin(x/\theta, n) \cdot Beta(\theta/a_D, b_D) \cdot d\theta$$

$$= \, ^nC_x \frac{1}{B(a_D, b_D)} \int \theta$$

$$P_n(x/D) = \int Bin_{n,\theta}(x) \cdot Beta_{a_D, b_D}(\theta) \cdot d\theta$$

$$= \, ^nC_x \cdot \frac{1}{B(a_D, b_D)} \cdot \int \theta^x (1-\theta)^{n-x} \cdot \theta^{a_D - 1} \cdot (1-\theta)^{b_D - 1} \cdot d\theta$$

$$= \, ^nC_x \cdot \frac{1}{B(a_D \cdot b_D)} \cdot B(x + a_D, n - x + b_D).$$

when $n = 1$ then

$$P(x = 1/D) = \frac{B(a_D + 1, b_D)}{B(a_D, b_D)} = \frac{\Gamma(a_D + 1)\Gamma(b_D)}{\Gamma(a_D + 1 + b_D)} \cdot \frac{\Gamma(a_D + b_D)}{\Gamma(a_D)\Gamma(b_D)}$$

$$= \frac{a_D}{a_D + b_D} \quad \left( \because \Gamma(a) = (a-1)\Gamma(a-1) \right)$$

$X \sim$ Number of Heads $\sim Bin(x), n = 5$

$$P(x < 3/\theta) = P(x = 0/\theta) + P(x = 1/\theta) + P(x = 2/\theta)$$

$$= \, ^5C_0 (1-\theta)^5 + \, ^5C_1 \theta(1-\theta)^4 + \, ^5C_2 \cdot \theta^2 (1-\theta)^3.$$

Posterior: $P(\theta/x < 3) \propto P(\theta) \cdot p(x < 3/\theta) P(\theta).$

$$P(\theta) = Beta_{(1,1)}(\theta) = \frac{\theta^{1-\theta} \cdot (1-\theta)^{1-1}}{B(1,1)} = 1 = U(0,1)$$

As prior is uniform, then

$$P(\theta/x < 3) \propto P(x < 3/\theta)$$

$$\propto \, ^5C_0 (1-\theta)^5 + \, ^5C_1 \theta(1-\theta)^4 + \, ^5C_2 \theta^2 (1-\theta)^3$$

3.7) a)) Non-Bayesian model is

$$Pois_\lambda(x) = e^{-\lambda} \cdot \frac{\lambda^x}{x!} \qquad \theta = (\lambda)$$

likelihood: $P(\not D/\theta) = \prod_{i=1}^{N} e^{-\theta} \cdot \frac{\theta^{x_i}}{x_i!}$

prior: $P(\theta) = Ga_{(a,b)}(\theta) = \frac{b^a}{\Gamma(a)} \cdot \theta^{a-1} \cdot e^{-\theta b}$   $a, b$ – hyper params

posterior $P(\theta/\not D) \propto P(\not D/\theta) \, P(\theta)$

$$\propto \prod_{i=1}^{N} e^{-\theta} \cdot \frac{\theta^{x_i}}{x_i!} \cdot \left(\theta^{a-1} \cdot e^{-\theta b}\right)$$

$$\propto \frac{e^{-N\theta} \, \theta^{\sum_i x_i}}{\prod_i x_i!} \cdot \theta^{a-1} \cdot e^{-\theta b}$$

$$\propto \theta^{\sum_i x_i + a - 1} \cdot e^{-\theta(N+b)}$$

posterior is also in the form of prior

$Ga(\theta) \propto \theta^{a-1} \cdot e^{-\theta b}$ with parameters.

$P(\theta/\not D) = Ga(a_D, b_D) \implies a_D = \sum_i x_i + a, \quad b_D = N + b.$

b) $E[P(\theta/\not D)] = E[Ga_{(a_D, b_D)}(\theta)]$

$$= \frac{a_D}{b_D} = \frac{\sum_i x_i + a}{N + b}$$

When $a \to 0, b \to 0$ this becomes

$$= \frac{\sum x_i}{N}$$

5.10) a) non-Bayesian model
$$p_\theta(x) = U(0,\theta)$$

• prior, $p(\theta) = Km^K \cdot \theta^{-(K+1)} \mathbb{1}_{\{\theta \geq m\}}$

hyper-params $(K,m)$

with non-informative prior $(K,m) = (0,0)$

$$p(\theta) \propto \frac{1}{\theta}$$

In first problem, we're given $A = \{100\}$, so

likelihood: $P(A/\theta) = \prod_{i=1}^{N} \frac{1}{\theta} \mathbb{1}_{\{x_i \in [0,\theta]\}}$

posterior: $P(\theta/A) \propto P(A/\theta) \, P(\theta)$

$$\propto \prod_{i=1}^{N} \frac{1}{\theta} \cdot \mathbb{1}_{\{x_i \in [0,\theta]\}} \cdot K \cdot m^K \cdot \theta^{-(K+1)} \cdot \mathbb{1}_{\{\theta \geq m\}}$$

$$\propto \frac{K \cdot m^K}{\theta^{N+K+1}} \left( \prod_{i=1}^{N} \mathbb{1}_{\{x_i \in [0,\theta]\}} \right) \mathbb{1}_{\{\theta \geq m\}}$$

$$\propto \frac{K \cdot m^K}{\theta^{N+K+1}} \cdot \mathbb{1}_{\{\theta \geq \max(A)\}}$$

If $\max(A) \leq m$ :

$$\propto \frac{(N+K) m^{N+K}}{\theta^{N+K+1}} \cdot \mathbb{1}_{\{\theta \geq m\}} = \text{pareto}(\theta)$$
$$(N+K, m)$$

if $\max(A) \geq m$ :

$$\propto \frac{(N+K) \max(A)^{N+K}}{\theta^{N+K+1}} \cdot \mathbb{1}_{\{\theta \geq \max(A)\}} = \text{pareto}(\theta)$$
$$(N+K, \max(A))$$

a) In this problem $\mathcal{D} = \{100\}$

$$mean(\mathcal{D}) = 100, \quad N = 1$$

posterior prior $p(\theta) = Pa(0,0) \quad k=0, \ m=0$

posterior $p(\theta/\mathcal{D}) =$ as $mean(\mathcal{D}) \geq m$

$$\propto \frac{(N+k) \cdot mean(\mathcal{D})^{N+k}}{\theta^{N+k+1}} \cdot 1 \quad \{\theta \geq mean(\mathcal{D})\}$$

$$\propto \frac{100}{\theta^2} \cdot 1 \quad \{\theta \geq 100\}$$

$$\propto pareto(\theta / 1, 100).$$

b) Mean of $pareto(\theta/k, m)$ if $k > 1$:

$$\frac{km}{k-1}$$

But here our $k_p = 1$, so mean doesn't exist

mode of $pareto(\theta/k, m) = m = 100$

Median of $pareto(\theta/k, m) = m \sqrt[k]{2}$

$$= 100 \cdot \sqrt{2} = 200$$

c) posterior predictive

$$P(x') \propto \int P(x'/\theta) \cdot P(\theta/\mathcal{D}) \, P(\theta)$$

$$\propto \int P(x'/\theta) \cdot pareto(\theta/k_{D1}, m_D)$$

As we have two types of posterior for $p(\theta/\mathcal{D})$

we get $P(\theta') = \begin{cases} \dfrac{k_D}{(N'+k_D)\, m_D^{N'}} & \text{if } \max(\theta') \leq m_D \\[4mm] \dfrac{k_D \cdot m_D^{k_D}}{(N'+k_D)\, \max(\theta')^{N'+k_D}} & \text{if } \max(\theta') > m_D \end{cases}$

we are asked to predikt for

$$P(x) = \frac{1}{2m} \mathbb{1}\{x \leq 100\}$$

$k_D = 1, \; m_D = 100$
$N' = 1, \; \max(\theta') = x$

$$+ \frac{100.1}{2x^{\nu}} \mathbb{1}\{x \geq 100\}$$

d) $P(x=100) = \dfrac{1}{200} \cdot \mathbb{1}\{100 \leq 200\} + \dfrac{100}{2 \times 10^4} \mathbb{1}\{100 > 100\}$

$= \dfrac{1}{200}$

$P(50 \leq x = 50)$

$= \dfrac{1}{200} \cdot \mathbb{1}\{50 \leq 100\} + 0 = \dfrac{1}{200}$

$P(x = 150) = \dfrac{100}{2 \times (150)^2} = \dfrac{1}{445}$

e) i) $\#$ our data is discrete, we can choose discrete probability distribution rather than (uniform) continous distribution

2) we took an un-informative prior, best which, didn't effect much on posterior. So, informative prior may help to get better results.

**3.11)** Bayesian Analysis of Exp...

$$P_\theta(n) = \theta e^{-\theta x} \quad \text{for} \quad X = \{x \geq 0\}$$
$$\theta > 0$$

**a)** MLE.

log-likelihood

$$\ell(\theta) = \sum_{i=1}^{N} \log \theta \cdot e^{-\theta x_i}$$

$$= \sum_{i=1}^{N} \left(\log \theta - \theta x_i\right) = N\log\theta - \theta \sum_{i=1}^{N} x_i$$

Set, $\dfrac{\partial \ell}{\partial \theta} = 0$ to get the $\arg\max_{\theta} \ell(\theta)$

$$\hat{\theta}_{MLE} = \arg\max_{\theta} \ell(\theta)$$

$$\Rightarrow \frac{\partial}{\partial \theta}\left(N\log\theta - \theta \sum_{i=1}^{N} x_i\right) = 0$$

$$\frac{N}{\theta} - \sum_{i=1}^{N} x_i = 0 \quad \Rightarrow \hat{\theta}_{MLE} \leq 1$$

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

**b)** After observing dataset $D = \{5, 6, 4\}$

$$\hat{\theta}_{MLE} = \frac{3}{5+6+4} = \frac{3}{15} = \frac{1}{5}$$

**c)** we assume prior $P(\theta)$ also follows exponential

$$P(\theta) = \text{Exp}_\lambda(\theta) = \lambda \cdot e^{-\lambda\theta}$$

and. with our knowledge $E[\theta] = \frac{1}{3}$.

we already know, $E[\theta] = \frac{1}{\lambda} = \frac{1}{3} \Rightarrow \lambda = 3$

d) posterior $p(\theta/\mathcal{D}) \propto p(\mathcal{D}/\theta) p(\theta)$

$$\propto \prod_{i=1}^{N} (\theta \cdot e^{-\theta x_i}) \cdot \lambda \cdot e^{-\lambda \theta}$$

$$\propto \theta^N \cdot e^{-\theta(\sum_{i=1}^{N} x_i + \lambda)}$$

So, posterior has form of Gamma with

params $(N+1, \sum_{i=1}^{N} x_i + \lambda)$    $(\because Ga_{(a,b)}(\theta) \propto \theta^{a-1} \cdot e^{-\theta b})$

e) If we write exponential has Gamma distribution then both likelihood and prior look simple

$(\because Exp_{\lambda}(\theta) = Ga_{(1,\lambda)}(\theta))$.

$$P(\mathcal{D}/\theta) = \prod_{i=1}^{N} \theta \cdot e^{-\theta x_i} = \theta^N \cdot e^{-\theta \sum_{i=1}^{N} x_i}$$

$$P(\theta) = \lambda \cdot e^{-\lambda \theta} \propto \theta^{1-1} \cdot e^{-\lambda \theta}$$

Both are in Gamma distribution form, so prior is conjugate to likelihood.

f) $E[\theta/\mathcal{D}]$ of Gamma distribution with

params $(N+1, \sum_{i=1}^{N} x_i + \lambda)$ is simply

$(\because E[Ga_{(a,b)}(\theta)] = a/b)$

$$= \frac{N+1}{\sum_{i=1}^{N} x_i + \lambda}$$

In denominator of posterior mean, its weighted combination of prior and likelihood mean. MLE.

g) $\hat{\theta}_{MLE} = \frac{1}{\frac{1}{N}\sum_{i=1}^{N} x_i} = \frac{1}{\sum_{i=1}^{N} x_i / N}$

$$E[\theta/\mathcal{D}] = \frac{1}{\frac{\sum_{i=1}^{N} x_i + \lambda}{N+1}}$$

As our dataset is small and prior is informative, we have to use posterior mean for $\theta$. in this example.

3.15) prior $p(\theta) = Beta(a, b)$.

and we are given $E[\theta] = m$ and $Var(\theta) = v$

As we already know expectation & variance for beta as

$$E[\theta] = \frac{a}{a+b} = m$$

$$Var(\theta) = \frac{ab}{(a+b)^2(a+b+1)} = v.$$

We can solve $a, b$ using these two equations.

by setting $m = 0.7$ and $v = 0.2^2 = 0.04$

$$\frac{a}{a+b} = 0.7 \quad \Rightarrow \quad a = 0.7a + 0.7b$$

$$a = \frac{7}{3}b.$$

$$\frac{a \cdot b}{(a+b)^2(a+b+1)} = 0.04 \Rightarrow \frac{\frac{7}{3}b \cdot b}{(\frac{7}{3}b+b)^2(\frac{7}{3}b+b+1)} = 0.04$$

$$\Rightarrow \frac{\frac{7}{3}b^2}{\frac{100}{9}b^2 \cdot (\frac{7}{3}b+b+1)} = 0.04$$

$$21 = 9(\frac{7}{3}b+b+1) \Rightarrow \frac{10b+3}{3} \cdot 9 = 21$$

$$10b+3 = \frac{63}{9}$$

$$10b = \frac{63}{9} - 3 = \frac{51}{9} = (2.75)$$

$$b = 1.275$$

$$a = \frac{7}{3} \cdot 1.275 = 2.975$$

4.14) $p(x) = N_\mu(\mu, \sigma^2)$ $\sigma^2$ - is given

and prior for $\mu$ is assumed as

$$p(\mu) = N(m, s^2)$$ hyper-params $(m, s^2)$.

From the derivation in class, posterior of $\mu$ is gaussian in form

$$p(\mu/\mathcal{D}) \propto e^{-\left(\mu - \frac{\frac{\sum x_i}{2\sigma^2} + \frac{m}{2s^2}}{\frac{n}{2\sigma^2} + \frac{1}{2s^2}}\right) \Big/ \frac{1}{\frac{n}{2\sigma^2} + \frac{1}{2s^2}}}$$

a) The MAP estimate for this

is $\mu_{MAP} = \dfrac{\frac{\sum x_i}{2\sigma^2} + \frac{m}{2s^2}}{\frac{n}{2\sigma^2} + \frac{1}{2s^2}} = \left(\frac{n\bar{x}}{\sigma^2} + \frac{m}{s^2}\right)\left(\frac{n}{\sigma^2} + \frac{1}{s^2}\right)^{-1}$

So, $\mu_{MAP}$ is weighted average of sample average $(\bar{x})$ and prior mean $(m)$

b) when $n$ increases, then $\frac{m}{s^2}$ in numerator and $\frac{1}{s^2}$ in denominator values diminishes.

So $\mu_{MAP} \triangleq \dfrac{\frac{\sum x_i}{2\sigma^2}}{\frac{n}{2\sigma^2}} = \dfrac{\sum x_i}{n} = \mu_{MLE}$

c) When prior variance increases then $\frac{1}{s^2} \to 0$
$(s^2)$

So, the expression becomes again

$$\mu_{MAP} \triangleq \frac{\sum x_i/2\sigma^2}{n/2\sigma^2} = \frac{\sum x_i}{n} = \mu_{MLE}$$

d) when prior variance decreases then $\frac{1}{s^2}$ becomes bigger. So $\mu_{MAP} = \dfrac{m/2s^2}{1/2s^2} = m$

* Discriminative model based on exponential family with sufficient statistics $\psi(\cdot)$ and parameters for $\psi_i(y)$ is $w_i^T \phi(x)$

$$P(y/x) = e^{\sum_i w_i^T \phi(x) \cdot \psi_i(y) - A(\sum_i w_i^T \phi(x))}$$

$$= e^{\sum_i \left( \sum_j w_{ij}^T \phi_j(x) \right) \cdot \psi_i(y) - A(\sum_i w_i^T \phi(x))}$$

Let call $\{ w^T \phi(x) = W$, so

$$A(W) = \log \left( \int e^{\sum_i \left( \sum_j w_{ij}^T \phi_j(x) \right) \psi_i(y)} \, dy \right)$$

In Generative model, we have sufficient statistics as $\phi(x) \otimes \psi(y)$, So

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

$$\alpha \, P(x, y)$$

$$\alpha \, e^{V^T \left( \phi(x) \otimes \psi(y) \right) - A(V)}$$

$$\alpha \, e^{\sum_{i,j} V_{ij} \phi_j(x) \psi_i(y) - A(V)}$$

$$\alpha \, e^{\sum_i \left( \sum_j V_{ij} \phi_j(x) \right) \psi_i(y) - A(V)}$$

Here $A(V) = \log \left( \int e^{\sum_i \left( \sum_j V_{ij} \, \phi_j(x) \right) \psi_i(y)} \, dy \right)$

As we can see both of them are in same form with only parameters differ. So, we can say even in generative discriminative model $\phi(\cdot)$ are

ficient statistics for inputs though they are
acting as a variables in canonical parameters of
inputs.