

# Bayesian Methods for ML

①

$X \& Y$  are  $\perp$ :

$$P(X, Y) = \underbrace{P(X)}_{\text{Joint}} \underbrace{P(Y)}_{\text{Marginal}}$$

## Conditional Prob

$$P(X|Y) = \frac{\underbrace{P(X, Y)}_{\text{Joint}}}{\underbrace{P(Y)}_{\text{Marginal}}} \quad \text{OR}$$

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

$$\Rightarrow P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

## Chain Rule

$$P(X, Y) = P(X|Y)P(Y)$$

$$P(X, Y, Z) = P(X|Y, Z)P(Y|Z)P(Z)$$

$$P(X_1, \dots, X_N) = \prod_{i=1}^N P(X_i | X_1 \dots X_{i-1})$$

$$\text{OR } P(\text{1st Cusp} | \text{and 2nd cusp}) = \frac{P(\text{and 2nd cusp})}{P(\text{1st Cusp})}$$

$$= \frac{P(\text{1st Cusp})}{P(\text{2nd cusp})}$$

## Sum Rule

$$P(X) = \int_{-\infty}^{\infty} p(x, \lambda) d\lambda$$

↑  
Marginalization

## Bayes theorem

$\theta$ -parameters

$X$ - observations

$$P(\theta|X) = \frac{P(X|\theta) \underbrace{P(\theta)}_{\text{Prior}}}{P(X) \underbrace{P(X|\theta)P(\theta)}_{\text{Evidence}}}$$

↑  
Posterior

(probability of params  
after we observe  
the data)

chain Rule

Joint

likelihood

Prior (prior knowledge  
about params)  
for instance  
some params  
are distributed  
around zero



# Bayesian Vs Frequentist

(3)

Frequentist

Bayesian

(See W1-L2)

## Classification

Training:

$$P(\theta | X_{tr}, y_{tr}) = \frac{P(y_{tr} | X_{tr}, \theta) P(\theta)}{P(y_{tr} | X_{tr})}$$

$X_{tr}, y_{tr}$   
↓  
 $\theta$

Prediction: (Marginalize " $\theta$ ")

$$P(y_{ts} | X_{ts}, X_{tr}, y_{tr}) = \int P(y_{ts} | X_{ts}, \theta) P(\theta | X_{tr}, y_{tr}) d\theta$$

In Bayesian Approach, prediction is a weighted average of output of our model for all possible values of parameters.

\* Bayes formula can also lead to regularisation

$$P(\theta | X) = \frac{P(X | \theta) P(\theta)}{P(X)} \quad P(\theta) = \text{prior \& can act as regularizer}$$

\* Bayesian methods are good for online learning likelihood

$$\text{For point } x_k, \quad P(\theta | x_k) = \frac{P(x_k | \theta) P(\theta)}{P(x_k)}$$

↑  
Posterior

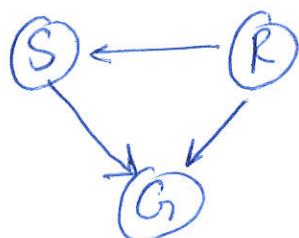
After computing posterior  $P(\theta|x_K)$ ; lets say we receive new point  $x_{K+1}$

$$\text{for } x_{K+1}, P(\theta|x_{K+1}) = \frac{P(x_{K+1}|\theta) P_K(\theta)}{P(x_{K+1})}$$

$P_K(\theta)$  can be taken as  $P(\theta|x_p)$  i.e. priors can be updated when new point arrives.

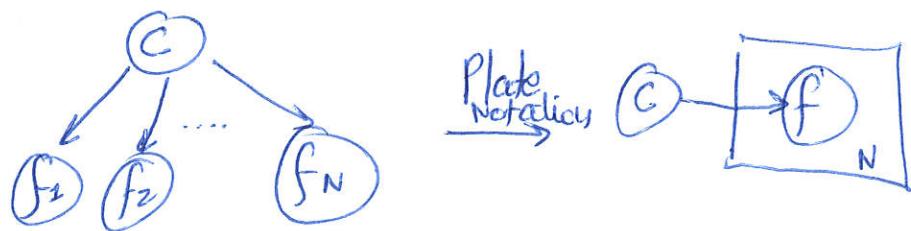
## How to define Model (Z-101)

Model: Joint Probability over all variables



$$P(S, R, G_1, G_2) = P(G_1|S, R) P(S|R) P(R)$$

## Naive Bayes Classifier



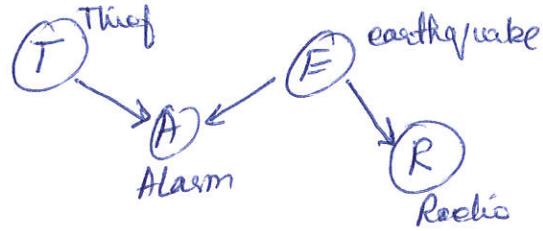
Joint Prob:

$$P(C, f_1, f_2, \dots, f_N) = P(C) \prod_{i=1}^N P(f_i|C)$$

- \* Bayesian Networks can not have/contain directed cycles. We can't have interdependent variables. For those cases we can either join random variables into one random variable or use HRF.

# W1-L3: Thief and Alarm

(5)



$$P(t, a, e, r) = P(t) P(e) P(a|t, e) P(r|e)$$

To fully define our model, we need to define these 4-probabilities  
i.e.  $P(t)$ ,  $P(e)$ ,  $P(a|t, e)$ ,  $P(r|e)$

Prior	
$P(T=1)$	$\frac{1}{1000}$
$P(E=1)$	$\frac{1}{100}$

P(a t, e)	
$P(A=1 \bar{T}, E)$	
No Thief: $\bar{T}=0$	$0 \quad 0.1$
Thief: $T=1$	$1 \quad 1$

$P(r e) \rightarrow$ radio report given earthq. wake	
$P(r e)$	
$E=0$	0
$E=1$	0.5

Calculate the following

\* i)  $P(T|A) = \text{thief given alarm}$

Using Bayes formula

$$P(T|A) = \frac{P(A|T) P(T)}{P(A)} = \frac{P(\bar{T}, A)}{P(A)}$$

Use sum rule

$$= \frac{P(\bar{T}, A, E) + P(\bar{T}, A, \bar{E})}{P(\bar{T}, A, E) + P(\bar{T}, A, \bar{E}) + P(\bar{T}, A, E) + P(\bar{T}, A, \bar{E})} \rightarrow \text{eq}(A)$$

In eq(A):

$$\begin{aligned} \text{i) } P(T, A, E) &= P(A|T, E) P(T) P(E) \\ &= (1) (10^{-3}) (10^{-2}) \\ \Rightarrow P(\bar{T}, A, E) &= 10^{-5} \end{aligned}$$

$$\begin{aligned} * P(\bar{T}, A, \bar{E}) &= P(A|\bar{T}, \bar{E}) P(\bar{T}) P(\bar{E}) \\ &= (1) (10^{-3}) (1 - 0.99) \\ &= 10^{-3} (0.99) \\ &= 0.00099 \end{aligned}$$

$$\begin{aligned} \text{ii) } P(\bar{T}, A, \bar{E}) &= P(A|\bar{T}, \bar{E}) P(\bar{T}) P(\bar{E}) \\ &= (0) ( ) ( ) \\ &= 0 \end{aligned}$$

$$\Rightarrow P(T|A) = \frac{P(\bar{T}, A, E) + P(\bar{T}, A, \bar{E})}{P(\bar{T}, A, E) + P(\bar{T}, A, \bar{E}) + P(\bar{T}, A, E) + P(\bar{T}, A, \bar{E})} \rightarrow 0 \approx 50\%$$

$$* P(T|A, R) = \text{ratio of joint Probabilities}$$

$$= \frac{P(A, T, R)}{P(A, R)} \quad \approx 1\%$$

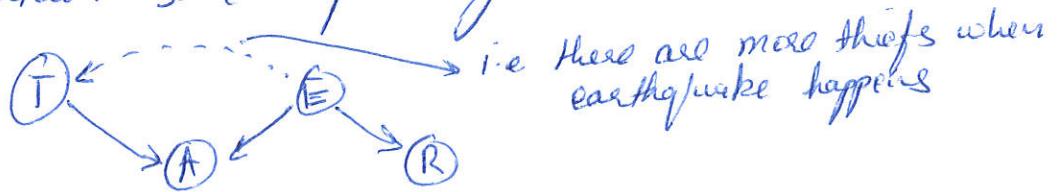
Add missing variables and using sum rule

$$= \frac{P(\bar{A}T, A, R, E) + P(T, A, R, \bar{E})}{P(A, R, \bar{T}, E) + P(A, R, T, \bar{E}) + P(A, \bar{R}, \bar{T}, E) + P(A, \bar{R}, T, \bar{E})} \rightarrow \text{Eq(B)}$$

In Eq(B):

$$\text{i) } P(T, A, R, E) = P(T) P(a|T, e) P(s|e) P(e)$$

1% for  $P(\text{thief}/\text{Alarm, radio repeat})$  suggests that our graphical model should contain some dependency like:

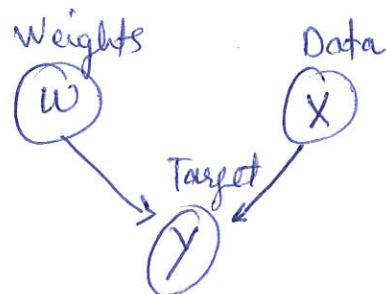


$$P(\bar{T}, A, R, E) = P(A|\bar{T}, E) P(R|E) P(T|E) P(E)$$

Least squares problem

$$L(w) = \sum_{i=1}^N (w^T x_i - y_i)^2 = \|w^T X - y\|^2 \rightarrow \min_w$$

$$\hat{w} = \operatorname{arg\,min}_w L(w)$$

Linear Regression in light of Bayesian N/WGMODEL

Three Random Variables

- i)  $w$
- ii)  $X$
- iii)  $y$

We are not interested in modelling the data ( $X$ ), therefore we can write the joint probability of weight & target given data as follows:

$$P(w, y | X) = P(y | X, w) P(w) \quad \left[ \text{See the model above} \right]$$

Note that  $P(w, y, X) = P(w, y | X)$  since we are not modelling data ( $X$ )

$$\text{In such case: } P(w, y, X) = P(y | w, X) P(w) P(X)$$

Since we don't model  $X$ ,

$$\Rightarrow P(w, y, X) = P(w, y | X) = P(y | w, X) P(w)$$

Now, we need to define  $P(y | w, X)$ . Let them be Normal distribution

$$\text{So } P(y | w, X) = \mathcal{N}(y | w^T X, \sigma^2)$$

The probability of target given weight & observation is gaussian centred at prediction  $w^T X$  & variance  $\sigma^2$

$$\text{and } P(w) = \mathcal{N}(w | 0, \gamma^2 I)$$

Now, let's train linear regression given the following formulae:

$$P(w, y|X) = P(y|X, w)p(w) \rightarrow \text{Eq(A)}$$

$$P(y|w, X) = N(y|w^T X, \sigma^2 I) \rightarrow \text{Eq(B)}$$

$$P(w) = N(w|0, \gamma^2 I) \rightarrow \text{Eq(C)}$$

Let's calculate the posterior probability of  $w|y, X$

$$P(w|y, X) = \frac{P(w, y, X)}{P(y, X)}$$

Since we are not modeling  $X$

$$P(w|y, X) = \frac{P(y, w|X)}{P(y|X)}$$

We want to find  $\max_w P(w|y, X)$

$$\Rightarrow w^* = \max_w \frac{P(y, w|X)}{P(y|X)} \quad \text{x when maximizing } \frac{P(y, w|X)}{P(y|X)}, \\ P(y, w|X) \text{ depends on } w.$$

In order to  $\max_w P(w|y, X)$ , we can maximize  $P(y, w|X)$ .

$$\Rightarrow w^* = \max_w P(y, w|X) \\ = \max_w P(y|X, w)p(w) \rightarrow \text{From eq(A)}$$

Applying log

$$\begin{aligned} w^* &= \max_w \log(P(y|X, w)p(w)) \\ &= \max_w \log P(y|X, w) + \log p(w) \\ &= \max_w \log C_1 \cdot \exp\left(-\frac{1}{2} (y - w^T X)^T [(\sigma^2 I)^{-1}] (y - w^T X)\right) \\ &\quad + \log C_2 \cdot \exp\left(-\frac{1}{2} w^T [\gamma^2 I]^{-1} w\right) \\ &= \max_w \left[ -\frac{1}{2\sigma^2} (y - w^T X)^T (y - w^T X) - \frac{1}{2\gamma^2} w^T w \right] \\ &= \max_w \left[ -\frac{1}{2\sigma^2} \|y - w^T X\|^2 - \frac{1}{2\gamma^2} \|w\|^2 \right] \end{aligned}$$

$$\begin{aligned} P(w|y, X) &= \frac{P(w, y, X)}{P(y, X)} \\ &= \frac{P(w|y|X)}{P(y|X)} \\ &= \dots \end{aligned}$$

Multiply by  $-1 \times 2\sigma^2$

$$= \min_{\omega} \underbrace{\|y - \omega^T x\|^2}_{\text{Sum of squares}} + \lambda \underbrace{\|\omega\|^2}_{L2-\text{regularizer}}$$

$$\Rightarrow \omega^* = \max_{\omega} p(\omega | y, x) = \min_{\omega} \left[ \|y - \omega^T x\|^2 + \lambda \|\omega\|^2 \right]$$

So, adding a normal prior on the weights [ $p(\omega)$ ], we turned from the least squares problem to the L2-regularized least-squares problem. i.e solving a least-square problem with L2-regularizer is equivalent of finding the MAP estimate for  $\omega$  with prior distribution  $\mathcal{N}(\omega | 0, \gamma I)$ .

## Conjugate Priors ~~(W1-L4)~~

### Analytical Inference (W1-L1)

\* Bayesian Inference has some complications. One of them is  $P(X)$  i.e computation of Evidence.

Posterior Distribution

$$P(\theta | X) = \frac{\text{Likelihood} \cdot \text{Prior}}{P(X) \text{ evidence}} = \frac{P(X|\theta) P(\theta)}{P(X)}$$

Computing this evidence is hard & complicated. For instance, in images if we model  $P(X)$  then we can draw new images. Generally computing  $P(X)$  is hard. We will see how can we avoid calculation of evidence during inference. It is called Maximum a posteriori principle.

### MAXIMUM A POSTERIORI

$$\theta_{MP} = \operatorname{argmax}_{\theta} P(\theta | X)$$

$$\theta_{MP} = \operatorname{argmax}_{\theta} \frac{P(X|\theta) P(\theta)}{P(X)}$$

$$\theta_{MP} = \operatorname{argmax}_{\theta} P(X|\theta) P(\theta) : \text{This is independent of } "x" \text{ so we can avoid computing } P(X)$$

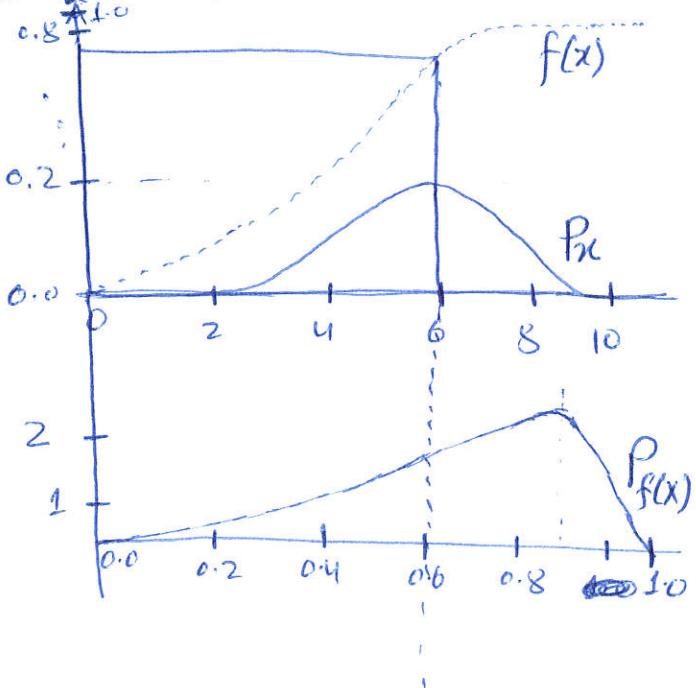
Thus it becomes an optimization problem.

$$\theta_{NP} = \underset{\theta}{\operatorname{argmax}} P(x|\theta) P(\theta)$$

### MAP PROBLEMS

Maximum a posteriori has its own problem. It is not invariant to reparametrization.

- i) Imagine we have a random variable "x" having Gaussian distribution as follows. we will apply sigmoid to it  $f(x)$  and get  $P_f(x)$  as an output. Note that the position of "MAXIMUM" changes ~~therefore we can't use it as prior distribution~~.



- ii) Next problem is we can't use it as a prior

$$P_K(\theta) = \frac{P(x_k|\theta) P_{K-1}(\theta)}{P(x_k)}$$

If we try to use MAP as prior, we will get the delta function and so we will not get any new information

$$P_K(\theta) = \frac{P(x_k|\theta) \delta(\theta - \theta_{NP})}{P(x_k)} = \delta(\theta - \theta_{NP})$$

Bayes formula:

$$P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \rightarrow \begin{array}{l} \text{Fixed by model} \\ \text{Our own choice} \\ \text{Fixed by data} \end{array}$$

lets select  $P(\theta)$  in a way that it becomes easier to compute the posterior:

$$\xrightarrow{\mathcal{A}(\nu')} P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \downarrow \mathcal{A}(\nu)$$

The prior is conjugate to the likelihood if the prior & the posterior lie in the same family of distributions. For instance if the prior is normal with some "mean" " $\theta$ " then we would expect the posterior to also be normal but with some other mean & variance.

Example

$$P(X|\theta) = N(X|\theta, \sigma^2)$$

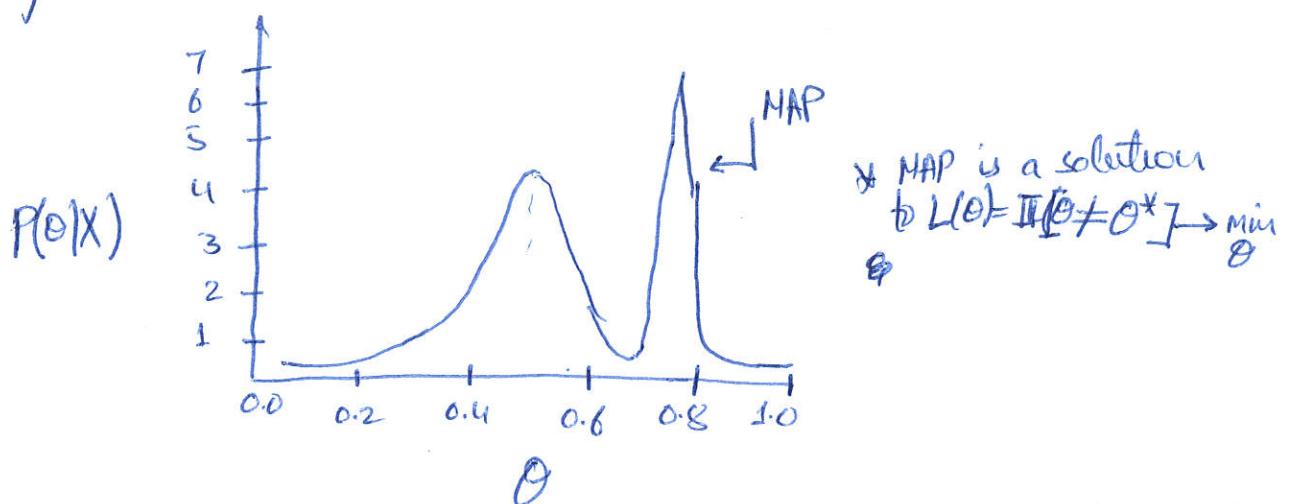
$$\mathcal{A}(\nu) = ??$$

$$\xrightarrow{\mathcal{A}(\nu')} P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} \rightarrow N(\theta, \sigma^2) \downarrow \mathcal{A}(\nu)$$

What should be the conjugate prior  $\mathcal{A}(\nu)$ ?

- A: If we select  $\mathcal{A}(\nu)$  as a normal distribution then  $P(\theta|X)$  will also have normal distribution. It is because product of two normal distributions  $P(X|\theta) \& P(\theta)$  yields normal distribution  $N(X|\theta, \sigma^2)$
- i.e.  $P(\theta|X) = \frac{P(X|\theta) P(\theta)}{f(X)} \leftarrow N(\theta|m, s^2)$
- $$\downarrow \mathcal{A}(\nu) \quad N(\theta|a, b^2)$$

iii) Another problem is that Maximum A posteriori estimation is actually an untypical point since there may be not enough probability density around it.



\* MAP is a solution  
to  $L(\theta) = \mathbb{I}[\theta \neq \theta^*] \rightarrow \min_{\theta}$

It also equals to the result of the minimization of the indicator that you do not end up in the true ideal parameter value.

We could use other functions for e.g. the squared/Absolute error. These would lead to the mean or median of the posterior distribution

### Objectives

$$L(\theta) = \mathbb{I}[\theta \neq \theta^*] \rightarrow \min_{\theta}$$

### Solution

Mode

$$L(\theta) = E[(\theta - \theta^*)^2] \rightarrow \min_{\theta}$$

Mean

$$L(\theta) = E|\theta - \theta^*| \rightarrow \min_{\theta}$$

Median

And if you choose these functions, we need to calculate evidence. That's what we are trying to avoid.

iv) finally we can't compute credible regions i.e. if  $D_{HP} = 12.53$  we can't say how confident are we.

## Summarizing Maximum A posteriori Estimation

### PROS

\* Easy to Compute

### CONS

- \* Not invariant to reparametrization
- \* Can't use as a prior
- \* finds untypical point
- \* Can't compute credible regions

\* Another way to avoid computing evidence: conjugate distributions

### Example

(13)

$$P(\theta|x) = \frac{p(x|\theta)p(\theta)}{P(x)} = \frac{N(x|\theta, 1)N(\theta|0, 1)}{P(x)}$$

$$p(\theta|x) \propto e^{-\frac{1}{2}(x-\theta)^2} e^{-\frac{1}{2}(\theta)^2}$$

$$p(\theta|x) \propto e^{-(\theta - \frac{x}{2})^2}$$

i.e  $p(\theta|x)$  is Normal with mean  $\frac{x}{2}$  & variance  $\frac{1}{2}$ . This also implies prior is conjugate to likelihood.

Question  $p(\theta)$  is Gaussian,  $p(x|\theta)$  is also Gaussian with mean  $\theta$  & variance  $\sigma^2$ . Here  $P(\theta)$  is not conjugate to  $P(X|\theta)$ .

W1-P2-L3Example: Normal, precision

\* More examples for conjugate priors

\* Gamma Distribution (parametrized by  $a, b$ )

Its PDF is:

$$\Gamma(y|a,b) = \frac{b^a}{\Gamma(a)} y^{a-1} e^{-by}$$

$y, a, b > 0$  (Gamma distribution is the distribution over the axis)

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

Statistics

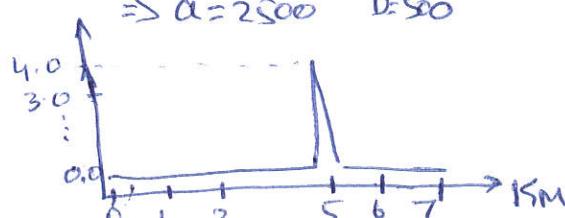
$$E[Y] = a/b \quad \underbrace{\text{mean}}_{\text{Mode}[Y] = \frac{a-1}{b}} \quad \text{Var}[Y] = a/b^2$$

Support of Gamma function distribution =  $[0, \infty)$ 

Example You run  $5\text{ km} \pm 100\text{ m}$  a day. (Model it by  $\Gamma(y|a, b)$ )

$$E[x] = a/b = 5 \quad \text{Var}[x] = a/b^2 = 0.1^2$$

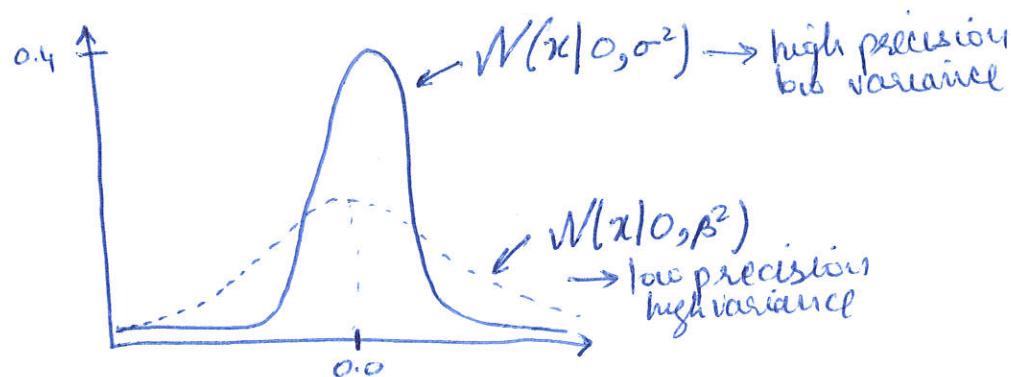
$$\Rightarrow a = 2500 \quad b = 500$$



Statement: The gamma distribution is conjugate to the normal w.r.t precision.

$$\text{Precision} = \frac{1}{\text{Variance}} = \frac{1}{\sigma^2}$$

$$\text{OR} \\ \gamma = \frac{1}{\sigma^2}$$



### Precision

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

$$\text{As } \sigma^2 = 1/y$$

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

Now, what is the conjugate prior w.r.t to the precision?

Ans: Here is functional form

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

Dropping all constants

$$N(x|\mu, y^{-1}) \propto \sqrt{y} e^{-\frac{y(x-\mu)^2}{2}}$$

$$\propto y^{1/2} e^{-by}$$

For final example, we will take Beta distribution

$$B(x|a,b) = \frac{1}{B(a,b)} x^{a-1} (1-x)^{b-1}$$

$$x \in [0,1]$$

$$a, b > 0$$

You can get multiple distributions like constant, U-shaped, uniform.

Statistics

$$E[x] = \frac{a}{a+b} \quad \text{Mode}[x] = \frac{a-1}{a+b-2}, \quad \text{Var}[x] = \frac{ab}{(a+b)^2(a+b-1)}$$

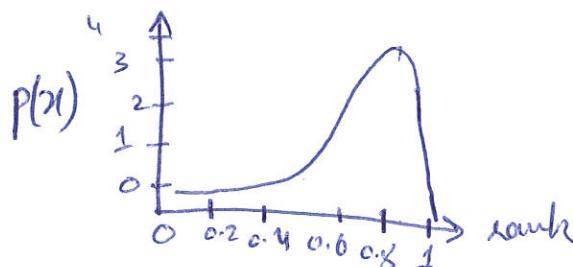
Support of beta distribution =  $[0,1]$

Note: Use Beta Distribution to model our favorite movie rank

Movie rank is  $0.8 \pm 0.1$

$$E[x] = \frac{a}{a+b} = 0.8$$

$$\text{Var}[x] = 0.1^2 \Rightarrow a=12, b=3$$



Actually, the Beta distribution is conjugate of the Bernoulli likelihood

$$p(x|\theta) = \theta^{N_1} (1-\theta)^{N_0} \rightarrow \text{Bernoulli likelihood} : \begin{matrix} N_1 \rightarrow \# \text{ of ones in dataset} \\ N_0 \rightarrow \# \text{ of zeros} \end{matrix}$$

lets select Beta distribution as prior

$$p(\theta) = B(\theta|a,b) \propto \theta^{a-1} (1-\theta)^{b-1}$$

Thus posterior :

$$p(\theta|x) \propto p(x|\theta) p(\theta)$$

$$\Rightarrow p(\theta|x) \propto \theta^{N_1} (1-\theta)^{N_0} \cdot \theta^{a-1} (1-\theta)^{b-1}$$

$$\Rightarrow p(\theta|x) \propto \theta^{N_1+a-1} (1-\theta)^{N_0+b-1}$$

This is Beta distribution with following params

$$p(\theta|x) = B(N_1+a, N_0+b)$$

lets try finding conjugate distribution in the following way:

$$p(y) \propto y^{\frac{1}{2}} e^{-by}$$

Use the Bayes formula :

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} \propto y e^{-y(b + \frac{(x-\mu)^2}{2})}$$

After dropping constants & rearranging terms

As we have seen that  $y$  has power " $\frac{1}{2}$ " instead of  $(\frac{1}{2})$  → this means it does not lie in the same family of distribution

Now, what if we choose gamma distributions. We will have another parameter that would allow us to vary the power of gamma

$$\text{As } p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

$$\text{where } N(x|\mu, y^{-1}) \neq \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}} \propto y^{\frac{1}{2}-by}$$

$$p(y) \propto y^{a-b-by}$$

This  $p(y)$  is gamma distribution with following params:

$$p(y) = \Gamma(y|a,b)$$

Now, lets try to compute our posterior  $p(y|x)$ :

$$\text{our prior is } p(y) : p(y) = \Gamma(y|a,b) \propto y^{a-1} e^{-by}$$

$$\text{our posterior } p(y|x) \propto p(x|y)p(y) \quad \begin{cases} \text{proportional to} \\ \text{likelihood \& prior} \end{cases}$$

Dropping all constants yields:

$$p(y|x) \propto \left( y^{\frac{1}{2}} e^{-y(b + \frac{(x-\mu)^2}{2})} \right) \cdot \left( y^{a-1} e^{-by} \right)$$

$$\Rightarrow p(y|x) \propto y^{\frac{1}{2}+a-1} e^{-y(b + \frac{(x-\mu)^2}{2})}$$

$$p(y|x) = \Gamma(a + \frac{1}{2}, b + \frac{(x-\mu)^2}{2})$$

mean      variance

So we avoid computing the evidence by choosing the conjugate prior and calculate posterior  $p(y|x)$  successfully.

## Overall Picture

(17)

$$P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)}$$

We would like to compute posterior  $P(\theta|X)$  but computing evidence  $P(X)$  makes it hard. Thus we choose conjugate priors that helped us computing  $P(\theta|X)$  without computing the evidence  $P(X)$ .

## Conjugate Pros & Cons

### PROS

- \* Exact posterior is computed using conjugate
- & It is easy for online-learning  
E.g  $P(\theta|X) = B(N_a + a, N_b + b)$

### CONS

- \* Conjugate Prior may be inadequate for some models.  
Therefore in the next week, we will see more advanced techniques to compute the full posterior or sometimes the approximate posterior.

