

$$P(x|\theta)$$

?

$$D = \{x_1, \dots, x_n\} \quad i.i.d$$

frequentist

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

Bayesian

$\theta$  is random var.

$$\theta \sim P(\theta)$$

$$D$$

$$\theta \sim P(\theta|D)$$

$$\hat{\theta}_{ML} = \arg \max_{\theta} P(D|\theta) = \arg \max_{\theta} P(x_1, \dots, x_n | \theta) \stackrel{i.i.d}{=} \arg \max_{\theta} \prod_{i=1}^n P(x_i | \theta)$$

✓  $\hat{\theta} \rightarrow \theta$

$n \rightarrow \infty$

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

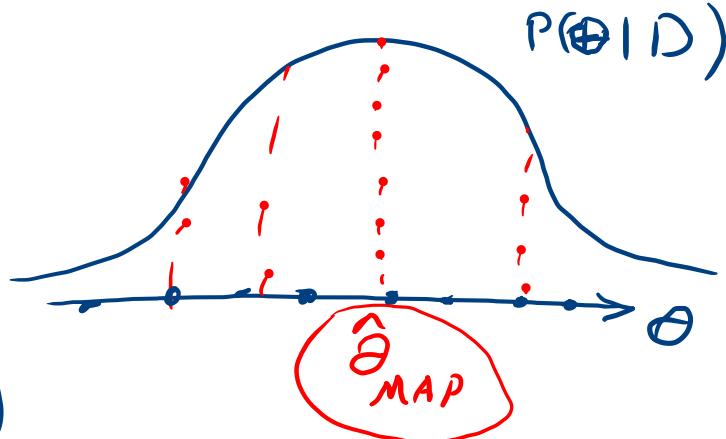
$$E[(\hat{\theta} - \theta)^2]$$

UMVE

$$\rightarrow P(x|\theta) = ?$$

$$P(\theta) \xrightarrow{D} \underline{P(\theta|D)}$$

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



$$\frac{P(x|D)}{\Downarrow} = \int P(x|\theta) P(\theta|D) d\theta = E_{\theta \sim P(\theta|D)} [P(x|\theta)]$$

predictive distribution

$$P(x|D) = \int \underbrace{P(x, \theta|D)}_{\downarrow} d\theta = \int \underbrace{P(x|\hat{\theta}, \hat{D})}_{\downarrow} \underbrace{P(\theta|D)}_{P(x|\theta)} d\theta$$

$$D \longrightarrow \theta \longrightarrow x$$

$$= \int P(x|\theta) P(\theta|D) d\theta$$

$$= E[P(x|\theta)]$$

$\Theta \sim P(\theta|D)$

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

1

$$\theta \sim P(\underline{\theta}) \xrightarrow{} P(\underline{\theta|D}) = \int f_i(\theta) P(\theta|D) d\theta = E[f_i(\theta)]$$

Conjugate Prior

$$\rightarrow \hat{\theta}_1$$

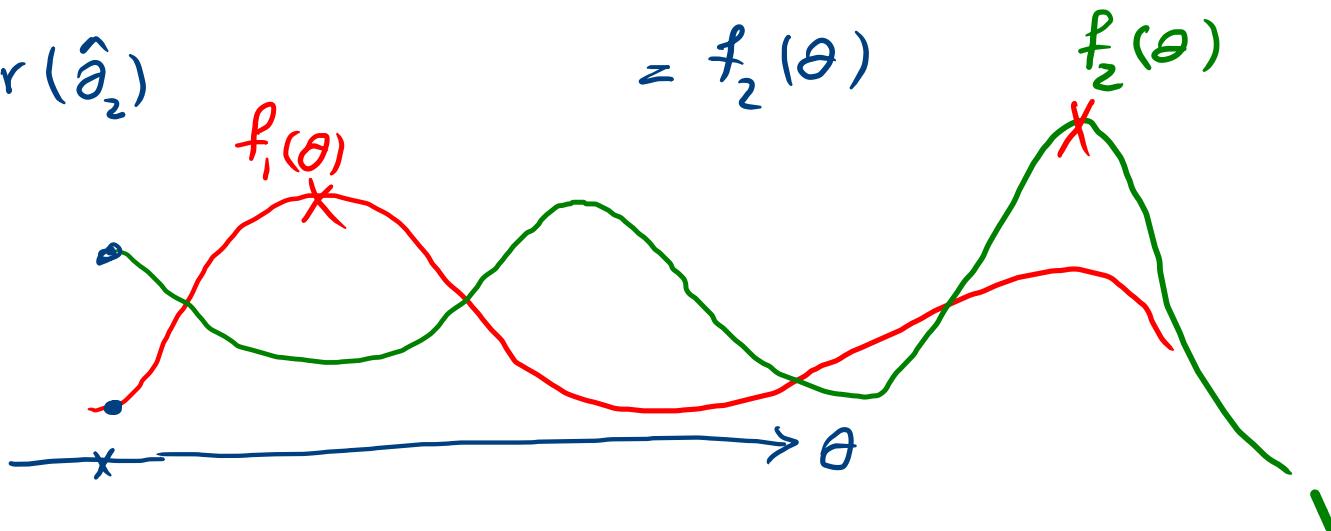
$$\text{var}(\hat{\theta}_1) = E[(\underline{\hat{\theta}} - \underline{\theta})^2] = f_1(\theta)$$

$$\rightarrow \hat{\theta}_2$$

$$\text{var}(\hat{\theta}_2)$$

$$= f_2(\theta)$$

minimax



$$P(x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!} \quad x \in \{0, 1, \dots\}$$

$$O\left(\frac{\alpha}{\beta^2}\right)$$

Gamma ( $\lambda | \alpha, \beta$ ) =  $c \lambda^{\alpha-1} e^{-\beta\lambda}$

hyper-parameter

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

$$\underline{P(D|\theta) P(\theta)} = \frac{n}{\prod_{i=1}^n} \underline{P(x_i|\theta)} P(\theta)$$

$$P(x_1, \dots, x_n | \theta)$$

$$= \frac{\lambda^{\sum x_i} e^{-n\lambda}}{\prod_{i=1}^n x_i!}$$

$$c \lambda^{\alpha-1} e^{-\beta\lambda}$$

$$= c' \lambda^{\sum x_i + \alpha - 1} e^{-(n+\beta)\lambda}$$

$n \rightarrow \infty$

$$= \text{Gamma}(\lambda \left| \sum_{i=1}^n x_i + \alpha, n + \beta \right.)$$

$$\alpha'$$

$$\beta'$$



$$\begin{aligned}
 p(\mu|\mathcal{D}) &= \frac{p(\mathcal{D}|\mu)p(\mu)}{\int p(\mathcal{D}|\mu)p(\mu) d\mu} \\
 &= \alpha \prod_{k=1}^n p(x_k|\mu)p(\mu),
 \end{aligned}$$

$\mathcal{N}(\underline{\mu}, \Sigma)$

$\mu \sim \mathcal{N}(\mu_0, \sigma_0^2)$

where  $\alpha$  is a **normalization factor** that depends on  $D$  but is independent of  $\mu$ .

This equation shows how the observation of a set of training samples **affects our ideas about the true value of  $\mu$** ;

$$\begin{aligned}
 p(\mu|\mathcal{D}) &= \alpha \prod_{k=1}^n \overbrace{\frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x_k - \mu}{\sigma} \right)^2 \right]}^{p(x_k|\mu) \sim \mathcal{N}(\mu, \sigma^2)} \overbrace{\frac{1}{\sqrt{2\pi}\sigma_0} \exp \left[ -\frac{1}{2} \left( \frac{\mu - \mu_0}{\sigma_0} \right)^2 \right]}^{p(\mu) \sim \mathcal{N}(\mu_0, \sigma_0^2)} \\
 &= \alpha' \exp \left[ -\frac{1}{2} \left( \sum_{k=1}^n \left( \frac{\mu - x_k}{\sigma} \right)^2 + \left( \frac{\mu - \mu_0}{\sigma_0} \right)^2 \right) \right] \\
 &= \alpha'' \exp \left[ -\frac{1}{2} \left[ \left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left( \frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \right) \mu \right] \right],
 \end{aligned}$$

factors that **do not depend on  $\mu$**  have been absorbed into the constants  $\alpha$ ,  $\alpha'$ , and  $\alpha''$ .



Given  $\mathcal{D} = \{x_1, \dots, x_n\}$ , we obtain

$$\begin{aligned} p(\mu|\mathcal{D}) &\propto \prod_{i=1}^n p(x_i|\mu)p(\mu) \\ &\propto \exp \left[ -\frac{1}{2} \left( \left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left( \frac{1}{\sigma^2} \sum_{i=1}^n x_i + \frac{\mu_0}{\sigma_0^2} \right) \mu \right) \right] \\ &= N(\mu_n, \sigma_n^2) \quad p(\mu) \text{ is said to be a } \textcolor{blue}{\textit{conjugate prior}} \end{aligned}$$

Where

$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \bar{x}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

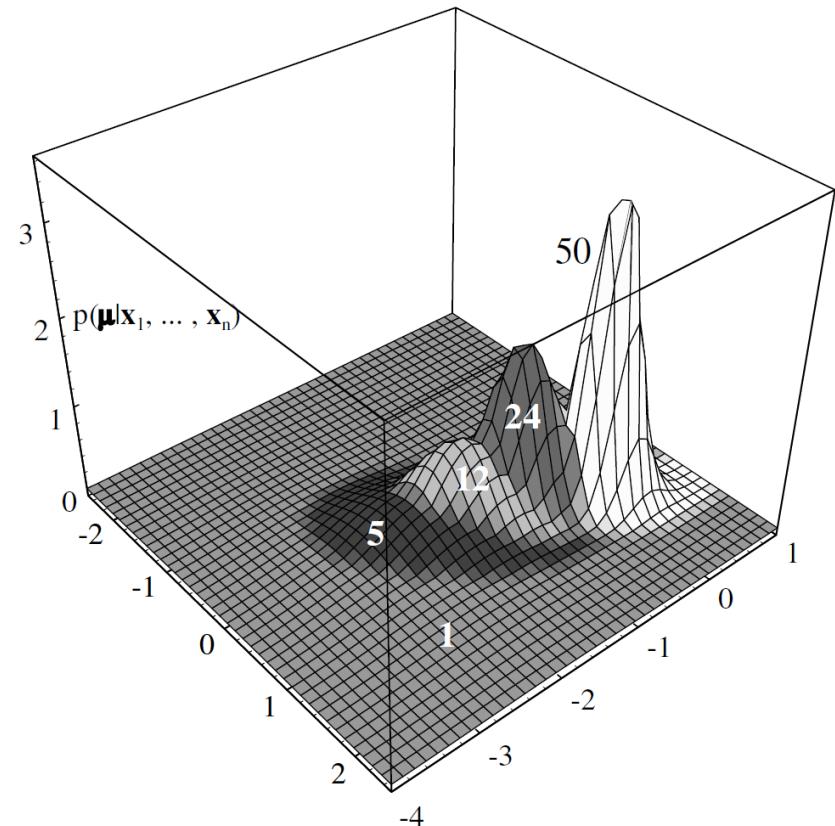
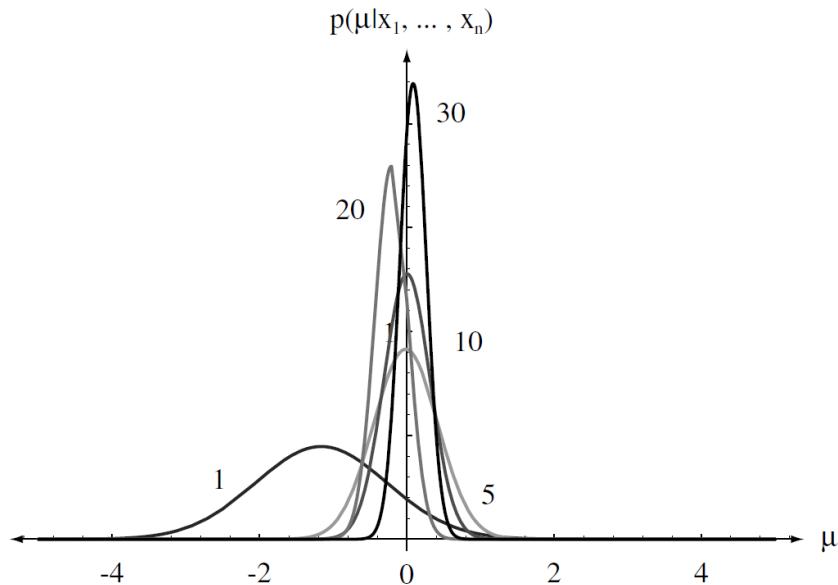
$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}.$$

# $\sigma_n^2$ uncertainty about $n^{\text{th}}$ guess

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n \sigma_0^2 + \sigma^2}.$$



- Since  $\sigma_n^2$  decreases monotonically with  $n$ - approaching  $\sigma^2/n$  as  $n$  approaches infinity
- Each additional observation decreases our uncertainty about the true value of  $\mu$ .
- As  $n$  increases,  $p(\mu | D)$  becomes more and more sharply peaked, approaching a Dirac delta function





# The Gaussian Case

- $\mu_0$  is our **best prior guess** and  $\sigma_0^2$  is the **uncertainty about this guess**.
- $\mu_n$  is our best guess **after observing n sample** in  $\mathbf{D}$  and  $\sigma_n^2$  is the uncertainty about this guess.
- $\mu_n$  always lies between  $\bar{x}_n$  and  $\mu_0$  with coefficients that are **non-negative and sum to one**.  
$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \bar{x}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$
  - If  $\sigma_0 = 0$ , then  $\mu_n = \mu_0$  (**no observation can change** our prior opinion).
  - If  $\sigma_0 \gg \sigma$ , then  $\mu_n = \bar{x}_n$  (we are **very uncertain** about our prior guess).
  - Otherwise,  $\mu_n$  approaches  $\bar{x}_n$  as n **approaches infinity**



# Class-conditional density

- Given the posterior density  $p(\mu|\mathcal{D})$ , the conditional density  $p(x|\mathcal{D})$  can be computed as

$$\begin{aligned} p(x|\mathcal{D}) &= \int p(x|\mu)p(\mu|\mathcal{D}) d\mu \\ &= \int \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2 \right] \frac{1}{\sqrt{2\pi}\sigma_n} \exp \left[ -\frac{1}{2} \left( \frac{\mu-\mu_n}{\sigma_n} \right)^2 \right] d\mu \\ &= \frac{1}{2\pi\sigma\sigma_n} \exp \left[ -\frac{1}{2} \frac{(x-\mu_n)^2}{\sigma^2 + \sigma_n^2} \right] f(\sigma, \sigma_n), \\ f(\sigma, \sigma_n) &= \int \exp \left[ -\frac{1}{2} \frac{\sigma^2 + \sigma_n^2}{\sigma^2\sigma_n^2} \left( \mu - \frac{\sigma_n^2 x + \sigma^2 \mu_n}{\sigma^2 + \sigma_n^2} \right)^2 \right] d\mu. \end{aligned}$$

predictive  
dist

- Where  $\rightarrow$  
$$p(x|\mathcal{D}) = N(\mu_n, \sigma^2 + \sigma_n^2)$$
- the **conditional mean  $\mu_n$**  is treated as if it were the true mean,
- the known **variance is increased** to account for our **lack of exact knowledge** of the mean  $\mu$ .



# The multivariate case

$$p(\mathbf{x}|\boldsymbol{\mu}) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

where  $\boldsymbol{\mu}$  is the only unknown parameter with a prior distribution

$$p(\boldsymbol{\mu}) = N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad (\boldsymbol{\Sigma}, \boldsymbol{\mu}_0 \text{ and } \boldsymbol{\Sigma}_0 \text{ are all known}).$$

Given  $D = \{x_1, \dots, x_n\}$ , we obtain

$$\begin{aligned} p(\boldsymbol{\mu}|D) \propto & \exp \left[ -\frac{1}{2} \left( \boldsymbol{\mu}^T \left( n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}_0^{-1} \right) \boldsymbol{\mu} \right. \right. \\ & \left. \left. - 2\boldsymbol{\mu}^T \left( \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n \mathbf{x}_i + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right) \right) \right]. \end{aligned}$$



# The Multivariate Gaussian Case

It follows that

$$p(\boldsymbol{\mu} | \mathcal{D}) = N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$$

where

$$\boldsymbol{\mu}_n = \boldsymbol{\Sigma}_0 \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \hat{\boldsymbol{\mu}}_n + \frac{1}{n} \boldsymbol{\Sigma} \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\mu}_0,$$

$$\boldsymbol{\Sigma}_n = \frac{1}{n} \boldsymbol{\Sigma}_0 \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\Sigma}.$$

$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \bar{x}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}.$$



# Class-conditional density

- Given the posterior density  $p(\mu | \mathcal{D})$ , the conditional density  $p(x | \mathcal{D})$  can be computed as
$$p(\mathbf{x} | \mathcal{D}) = N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma} + \boldsymbol{\Sigma}_n)$$
- Which can be viewed as the **sum of** a random vector  $\mu$  with
$$p(\boldsymbol{\mu} | \mathcal{D}) = N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$$
- and an **independent random vector  $y$**  with
$$p(\mathbf{y}) = N(0, \boldsymbol{\Sigma}).$$



# The Bernoulli Case

- Consider  $P(x|q) = \text{Bernoulli}(q)$  where  $q$  is the unknown parameter with a prior distribution

$p(q) = \text{Beta}(a, b)$  ( $a$  and  $b$  are both known).

$$f(x; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

- Given  $D = \{x_1, \dots, x_n\}$ , we obtain

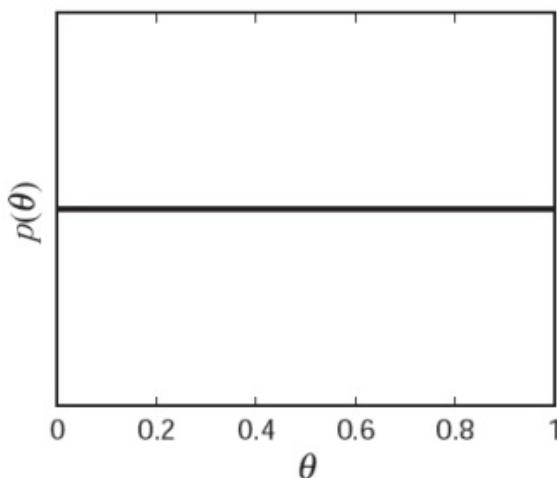
$$p(\theta|D) = \text{Beta}\left(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i\right).$$

$$x \sim \text{Bernoulli}(\theta) \Rightarrow P(x) = \underline{\theta}^x \ (1-\theta)^{1-x} \quad x \in \{0, 1\}$$

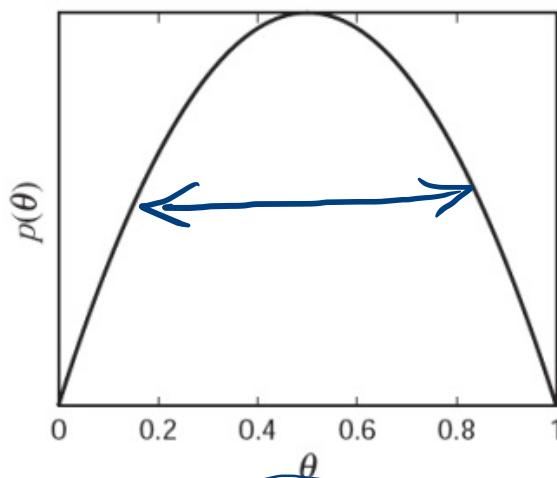
$$\theta \sim \underline{\underline{P(\theta)}} \quad 0 \leq \theta \leq 1$$

$$\text{Beta}(\alpha, \beta)$$

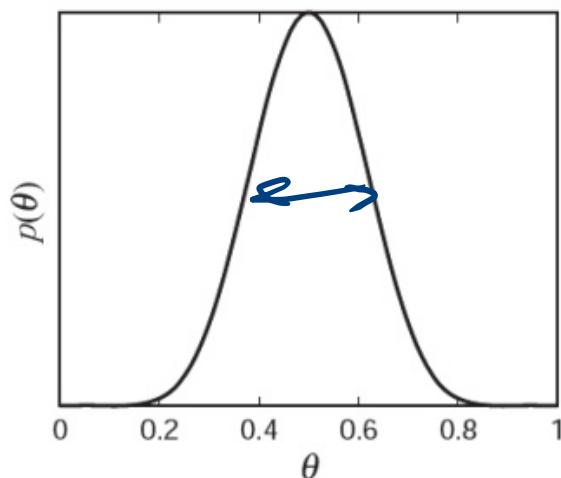
$\text{Beta}(\alpha, \beta)$



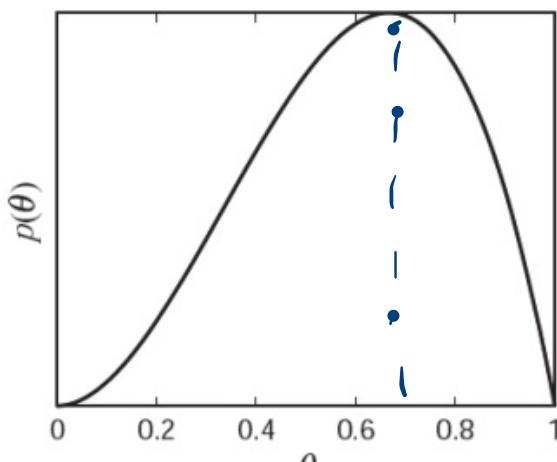
$\text{Beta}(1,1)$



$\text{Beta}(2,2)$

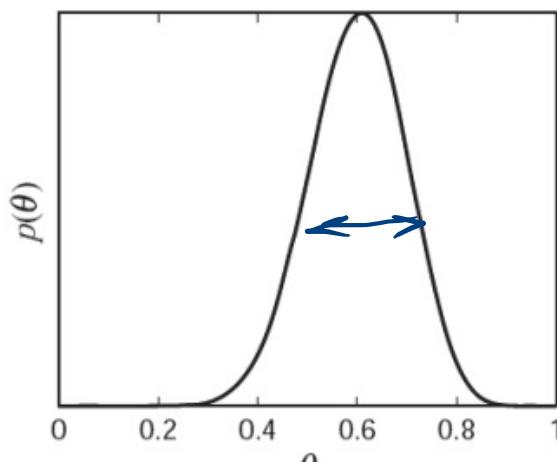


$\text{Beta}(10,10)$

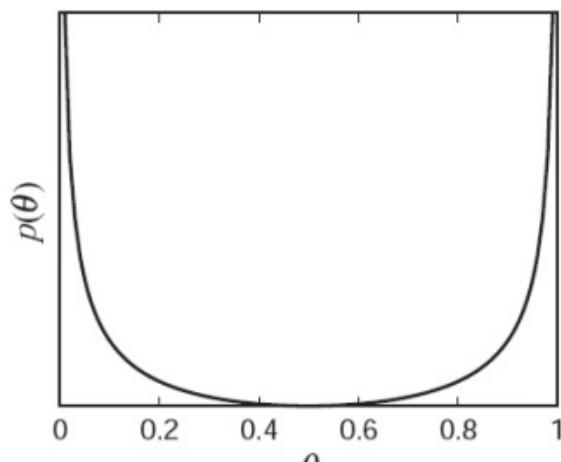


$\text{Beta}(3,2)$

$\alpha > \beta$

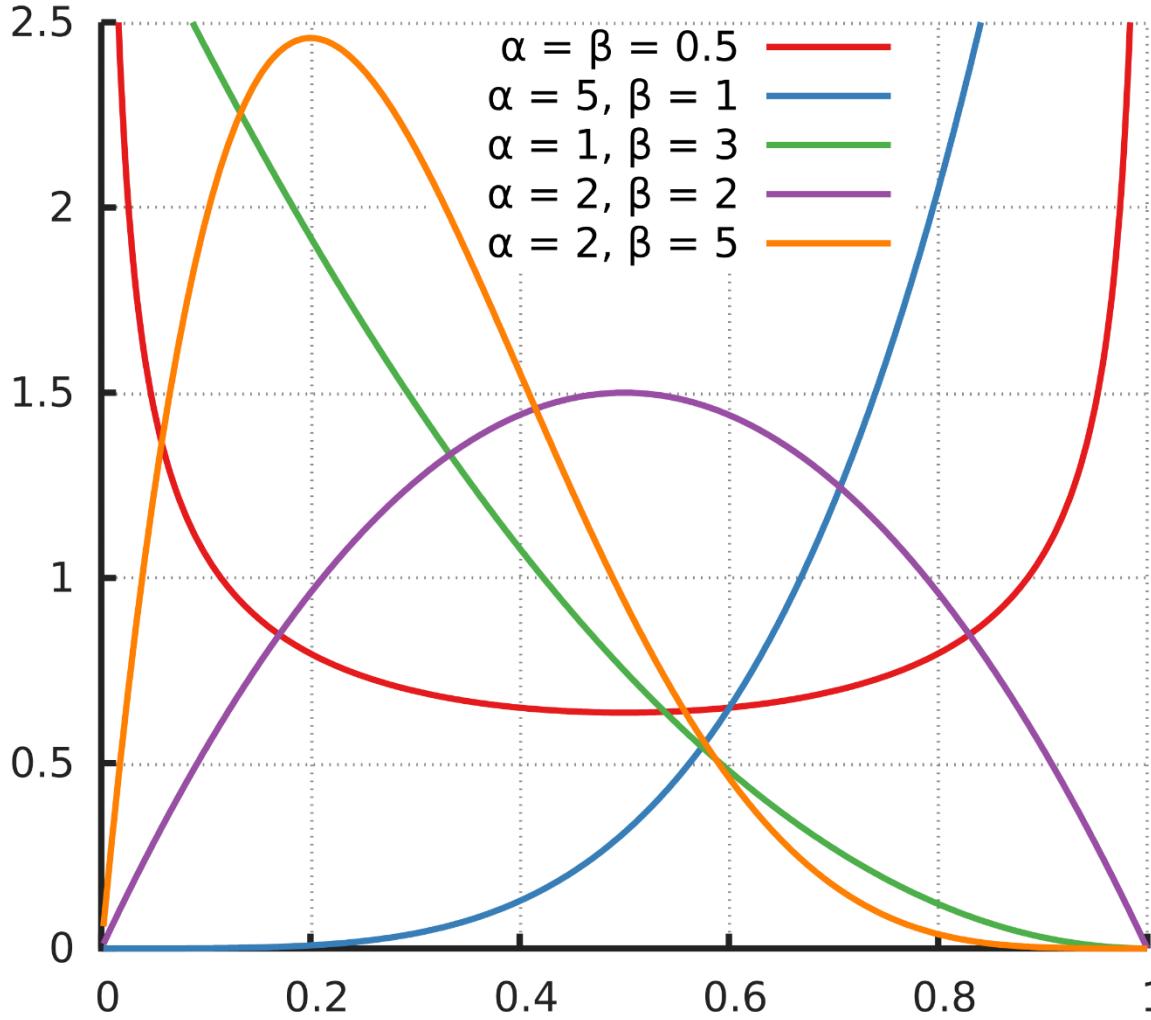


$\text{Beta}(15,10)$



$\text{Beta}(0.5,0.5)$

.



$$\frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$$

$\alpha > 0$  shape (real)

$\beta > 0$  shape (real)

$x \in [0, 1]$  or  $x \in (0, 1)$

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx,$$

$$E[X] = \frac{\alpha}{\alpha + \beta}$$

### Mode

$$\left| \begin{array}{l} \frac{\alpha - 1}{\alpha + \beta - 2} \text{ for } \alpha, \beta > 1 \\ \end{array} \right.$$

$$\text{var}[X] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$



# The Bernoulli Case

- The Bayes estimate of  $\mathbf{q}$  can be computed as the expected value of  $p(\mathbf{q} | D)$ , i.e.,

$$\begin{aligned}\hat{\theta} &= \frac{\alpha + \sum_{i=1}^n x_i}{\alpha + \beta + n} \\ &= \left( \frac{n}{\alpha + \beta + n} \right) \frac{1}{n} \sum_{i=1}^n x_i + \left( \frac{\alpha + \beta}{\alpha + \beta + n} \right) \frac{\alpha}{\alpha + \beta}.\end{aligned}$$

# Conjugate Priors

$$\mu \sim N(\dots)$$

$$(\sigma^2) \sim G(\dots)$$



- A conjugate prior is one which, when multiplied with the probability of the observation, gives a posterior probability **having the same functional form** as the prior.
- This relationship **allows the posterior to be used as a prior** in further computations.

$$P(\mu, \sigma^2)$$

$$= P(\mu)$$

$$P(\sigma^2)$$

<i>pdf generating the sample</i>	<i>corresponding conjugate prior</i>
Gaussian	Gaussian
Exponential	Gamma
Poisson	Gamma
Binomial	Beta
Multinomial	Dirichlet



# Recursive Bayes Learning

- What about the **convergence** of  $p(x|D)$  to  $p(x)$ ?
- Given  $\mathcal{D}^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , for  $n > 1$

$$p(\mathcal{D}^n|\boldsymbol{\theta}) = p(\mathbf{x}_n|\boldsymbol{\theta})p(\mathcal{D}^{n-1}|\boldsymbol{\theta})$$

and

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}) d\boldsymbol{\theta}},$$

$$p(\boldsymbol{\theta}|\mathcal{D}^n) = \frac{p(\mathbf{x}_n|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}^{n-1})}{\int p(\mathbf{x}_n|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}^{n-1}) d\boldsymbol{\theta}}$$

where

$$p(\boldsymbol{\theta}|\mathcal{D}^0) = p(\boldsymbol{\theta})$$

- Quite useful if the distributions can be represented using **only a few parameters (sufficient statistics)**.



# Recursive Bayes Learning

- Consider the Bernoulli case  $P(x|q) = \text{Bernoulli}(q)$  where  $p(q) = \text{Beta}(a, b)$ , the Bayes estimate of  $q$  is

$$\hat{\theta} = \frac{\alpha}{\alpha + \beta}.$$

- Given the training set  $D = \{x_1, \dots, x_n\}$ , we obtain

$$p(\theta|D) = \text{Beta}(\alpha + m, \beta + n - m)$$

where

$$m = \sum_{i=1}^n x_i = \#\{x_i | x_i = 1, x_i \in D\}.$$



# Recursive Bayes Learning

- The Bayes estimate of  $\theta$  becomes

$$\hat{\theta} = \frac{\alpha + m}{\alpha + \beta + n}.$$

- Then, given a **new training** set

$$\mathcal{D}' = \{x_1, \dots, x_{n'}\}$$

- We obtain

$$p(\theta|\mathcal{D}, \mathcal{D}') = \text{Beta}(\alpha + m + m', \beta + n - m + n' - m')$$

- Where

$$m' = \sum_{i=1}^{n'} x_i = \#\{x_i | x_i = 1, x_i \in \mathcal{D}'\}.$$



# Recursive Bayes Learning

- The Bayes estimate of  $q$  becomes

$$\hat{\theta} = \frac{\alpha + m + m'}{\alpha + \beta + n + n'}.$$

- Thus, recursive Bayes learning involves **only keeping the counts  $m$**  (related to **sufficient statistics** of Beta) and the number of training samples  $n$ .

# Comparison of MLEs and Bayes estimates



	<i>MLE</i>	<i>Bayes</i>
<i>computational complexity</i>	differential calculus, gradient search	multidimensional integration
<i>interpretability</i>	point estimate	weighted average of models
<i>prior information</i>	assume the parametric model $p(\mathbf{x} \boldsymbol{\theta})$	assume the models $p(\boldsymbol{\theta})$ and $p(\mathbf{x} \boldsymbol{\theta})$ but the resulting distri- bution $p(\mathbf{x} \mathcal{D})$ may not have the same form as $p(\mathbf{x} \boldsymbol{\theta})$

If there is **much data** (strongly peaked  $p(\mathbf{q}|\mathcal{D})$ ) and the prior  $p(\mathbf{q})$  is **uniform**, then the Bayes estimate and MLE are **equivalent**.



# Classification Error

- To apply these results to **multiple classes**, separate the training samples to  $c$  subsets  $D_1, \dots, D_c$ , with the samples in  $D_i$  belonging to class  $w_i$ , and then **estimate each density**  $p(x|w_i, D_i)$  **separately**.
- Different **sources** of error:
  - • **Bayes error**: due to overlapping class-conditional densities (related to the features used; inherent property of the problem and **can never be eliminated**).
  - • **Model error**: due to incorrect model.
  - • **Estimation error**: due to estimation from a finite sample (can be reduced by increasing the amount of training data).

$$\boxed{P(y|\alpha) \quad P(x|y) \quad P(y)}$$

# Logistic Regression:

{ Discriminative classifier  
Bayes classifier  
Linear classifier

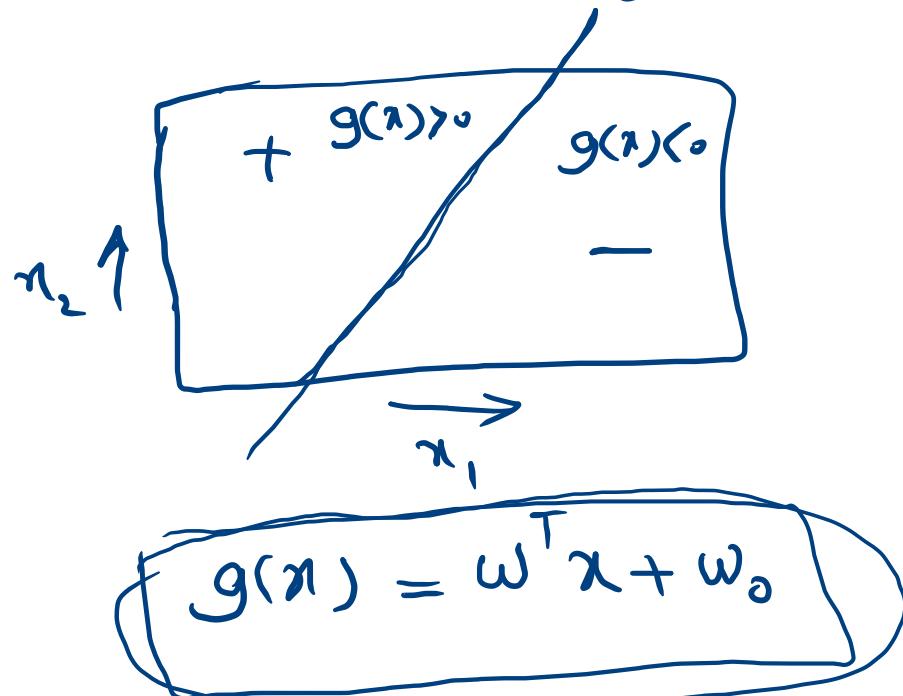
$$\frac{P(y=1|\alpha)}{P(y=0|\alpha)} \stackrel{!}{\geq} 1$$

$\checkmark$

$g(\alpha) \geq 0$

Binary  $\alpha$

$$g(\alpha) = \ln \frac{P(y=1|\alpha)}{1 - P(y=1|\alpha)} \stackrel{!}{\geq} 0$$



$$g(\alpha) = \omega^T \alpha + \omega_0$$

$$g(x) = \ln \frac{P(y=1|x)}{1 - P(y=1|x)} = \underline{\omega^T x}$$

$$P(y=1|x) = \frac{1}{1 + e^{-\omega^T x}}$$

$$(P(y=1|x)) = \sigma(\underline{\omega^T x})$$

✓

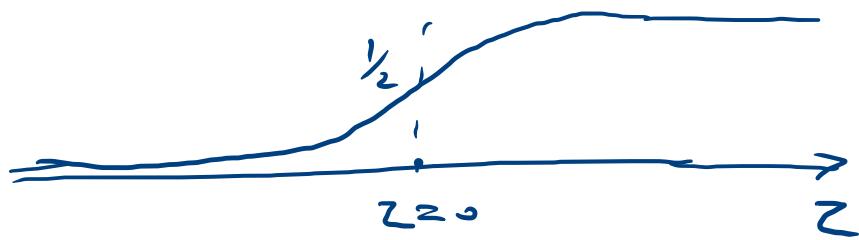
$$P(y|x) = \underline{\theta}^y (1-\underline{\theta})^{1-y} = \sigma(\underline{\omega^T x})^y (1 - \sigma(\underline{\omega^T x}))^{1-y}$$

$$g(x) = \omega^T x + \omega_0$$

$$\omega = \begin{bmatrix} \omega_0 \\ \omega_1 \\ \vdots \\ \omega_d \end{bmatrix}$$

$$x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



$$\cdot \left. \left( p(y|x) = \sigma(\omega^T x)^y (1 - \sigma(\omega^T x))^{1-y} \right) \right]$$

$$\rightarrow D = \{ (x_1, y_1), \dots, (x_n, y_n) \}$$

$$\hat{\omega}_{ML} = \arg \max_{\omega} \ln P(D|\omega)$$

$$\log P(D|\omega) = \log P(y_1, \dots, y_n | \underbrace{x_1, \dots, x_n}_{\text{i.i.d}}, \omega) \stackrel{\text{i.i.d}}{=} \sum_{i=1}^n \log \underline{P(y_i | x_i, \omega)}$$

$$= \sum_{i=1}^n y_i \log \sigma(\omega^T x_i) + (1 - y_i) \log (1 - \sigma(\omega^T x_i)) = L(\omega)$$

$$\nabla_{\omega} L(\omega) = \sum_{i=1}^n x_i \left( y_i - \frac{\sigma(\omega^T x_i)}{1 + e^{-\omega^T x_i}} \right) \quad \hat{\omega}_{ML}$$

$$P(\pi|\theta)$$

$$P(y|x, \theta)$$

$$\arg \max_{\theta} P(x_1, \dots, x_n | \theta)$$

$$P(y_1, \dots, y_n | x_1, \dots, x_n, \theta)$$