# **Machine Learning**
## 4771

Instructors:
Adrian Weller and Ilia Vovsha

# Topic 2: Basic concepts of Bayesians and Frequentists

- Properties of PDFs

- Bayesians & Frequentists

- ML, MAP and Full Bayes

- Example: Coin Toss

- Bernoulli Priors

- Conjugate Priors

- Bayesian decision theory

# Properties of PDFs

<span style="color:red">Probability Distribution Function</span>

- Review some basics of probability theory

- First, pdf is a function, multiple inputs, one output:

$$p\left(x_1, \ldots, x_n\right) \qquad p\left(X_1 = 0.3, \ldots, X_n = 1\right) = 0.2$$

- Function's output is always non-negative:

$$p\left(x_1, \ldots, x_n\right) \geq 0$$

- Can have discrete or continuous inputs or both:

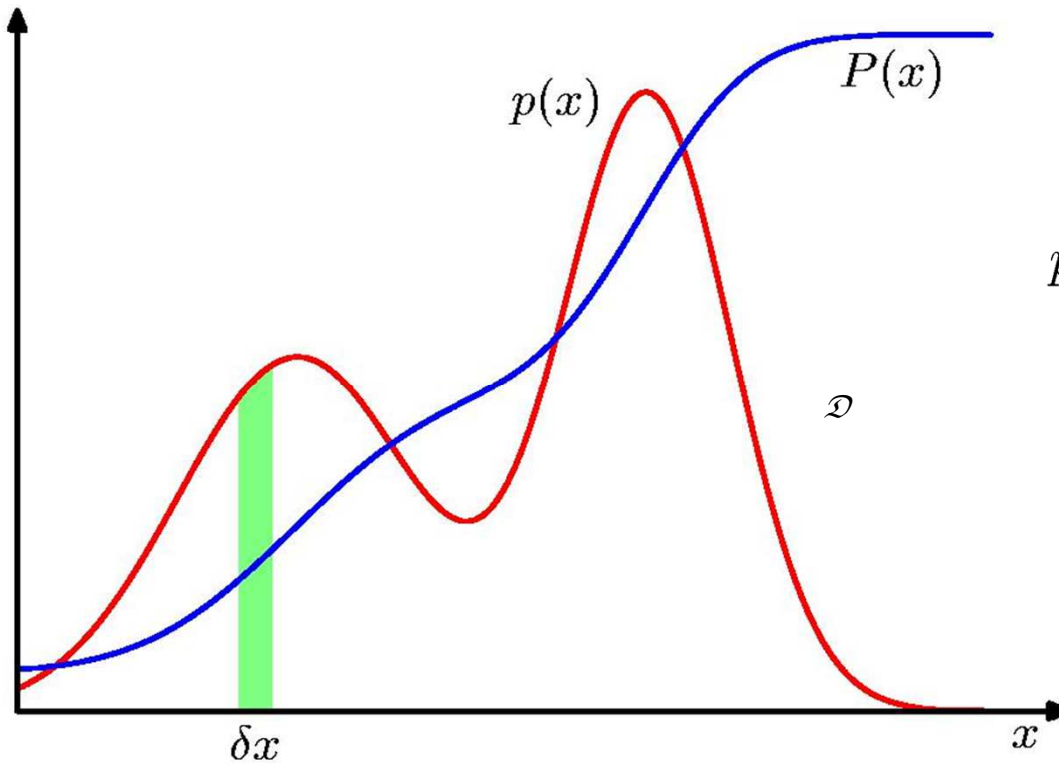$$p\left(X_1 = 1, X_2 = 0, X_3 = 0, X_4 = 3.1415\right)$$

- Summing over the domain of all inputs gives unity:

$$\int_{y=-\infty}^{\infty} \int_{x=-\infty}^{\infty} p\left(x, y\right) dx dy = 1 \qquad \sum_y \sum_x p\left(x, y\right) = 1$$

|   | | |
|---|---|---|
| Y 0 | 0.4 | 0.1 |
| 1 | 0.3 | 0.2 |

0       1

X

**Continuous→integral, Discrete→sum**

# Properties of PDFs



PDF

$$p(x \in (a, b)) = \int_a^b p(x)\, \mathrm{d}x$$

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

CDF

Cumulative
Distribution
Function

$$p(x) \geqslant 0 \qquad \int_{-\infty}^{\infty} p(x)\, \mathrm{d}x = 1$$

| Y | 0 | 0.4 | 0.1 |
|---|---|-----|-----|
|   | 1 | 0.3 | 0.2 |
|   |   | 0   | 1   |
|   |   |     | X   |

# Properties of PDFs

- Marginalizing: integrate/sum out a variable leaves a marginal distribution over the remaining ones…

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

- Conditioning: if a variable 'y' is 'given' we get a conditional distribution over the remaining ones…

$$p(x \mid y) = \frac{p(x,y)}{p(y)}$$

- Bayes Rule: mathematically just redo conditioning but has a deeper meaning (1764)… if we have x being data and $\theta$ being a model

**likelihood**

**posterior** $\longrightarrow$ $$p(\theta \mid x) = \frac{p(x \mid \theta) p(\theta)}{p(x)}$$ $\longleftarrow$ **prior**

**evidence**

# Properties of PDFs

- Expectation: can use pdf p(x) to compute averages and expected values for quantities, denoted by:

$$E_{p(x)}\left\{f(x)\right\} = \int_x p(x)f(x)dx \quad or \quad = \sum_x p(x)f(x)$$

- Properties:

$$E\left\{cf(x)\right\} = cE\left\{f(x)\right\}$$

$$E\left\{f(x)+c\right\} = E\left\{f(x)\right\}+c$$

$$E\left\{E\left\{f(x)\right\}\right\} = E\left\{f(x)\right\}$$

- Mean: expected value for x

$$E_{p(x)}\left\{x\right\} = \int_{-\infty}^{\infty} p(x)x\,dx$$

**example: speeding ticket**

| Fine=0$ | Fine=20$ |
|---------|----------|
| 0.8     | 0.2      |

**expected cost of speeding?**

- Variance: expected value of (x-mean)$^2$, how much x varies

$$Var\left\{x\right\} = E\left\{\left(x - E\left\{x\right\}\right)^2\right\} = E\left\{x^2 - 2xE\left\{x\right\} + E\left\{x\right\}^2\right\}$$

$$= E\left\{x^2\right\} - 2E\left\{x\right\}E\left\{x\right\} + E\left\{x\right\}^2 = E\left\{x^2\right\} - E\left\{x\right\}^2$$

# The IID Assumption

- Most of the time, we will assume that a dataset is independent and identically distributed (IID)

- In many real situations, data is generated by some black box phenomenon in an arbitrary order.
- Assume we are given a dataset:

$$\mathcal{D} = \left\{ x_1, \ldots, x_N \right\}$$

"Independent" means that (given the model $\theta$) the probability of our data multiplies:

$$p\left(x_1, \ldots, x_N \mid \Theta\right) = \prod_{i=1}^{N} p_i\left(x_i \mid \Theta\right)$$

"Identically distributed" means that each marginal probability is the same for each data point

$$p\left(x_1, \ldots, x_N \mid \Theta\right) = \prod_{i=1}^{N} p_i\left(x_i \mid \Theta\right) = \prod_{i=1}^{N} p\left(x_i \mid \Theta\right)$$

# Ex: Is a coin fair?

A stranger tells you his coin is fair.

Let's assume tosses are iid with P(H)=$\mu$.

He tosses it 4 times, gets H H T H.

What can you say about $\mu$?

# Bayesians & Frequentists

• Frequentists (Neymann/Pearson/Wald). An orthodox view that sampling is infinite and decision rules can be sharp.

• Bayesians (Bayes/Laplace/de Finetti). Unknown quantities are treated probabilistically and the state of the world can always be updated.

*actuarial fair*

de Finetti: p( event ) = price I would pay for a contract that pays $1 when event happens

• Likelihoodists (Fisher). Single sample inference based on maximizing the likelihood function.

# Bayesians & Frequentists

- Frequentists:
    - Data are a repeatable random sample - there is a frequency
    - Underlying parameters remain constant during this repeatable process
    - Parameters are fixed

- Bayesians:
    - Data are observed from the realized sample.
    - Parameters are unknown and described probabilistically
    - Data are fixed

# Frequentists

- Frequentists:  classical / objective view / no priors
  every statistician should compute same p(x) so no priors
  can't have a p(event) if it never happened
  avoid p($\theta$), there is 1 true model, not distribution of them
  permitted: $p_\theta(x,y)$   forbidden: $p(x,y|\theta)$
  Frequentist inference: estimate one best model $\theta$
      use the Maximum Likelihood Estimator (ML)
      (unbiased & minimum variance)
      do not depend on Bayes rule for learning

$$\theta_{ML} = \arg\max_\theta p(\mathcal{D}|\theta)$$

$$\text{Data } \mathcal{D} = (x_1, x_2, \ldots, x_n)$$

# Bayesians

- **Bayesians:** subjective view / priors are ok
  put a distribution or pdf on all variables in the problem
  even models & deterministic quantities (speed of light)
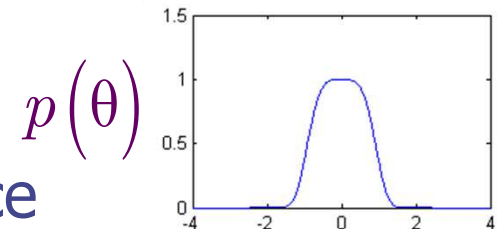  use a prior $p(\theta)$ on the model $\theta$ before seeing any data

  Bayesian inference: use Bayes rule for learning, integrate
  over all model $(\theta)$ unknown variables

# Bayesian Inference

- Bayes rule can lead us to maximum likelihood
- Assume we have a prior over models p($\theta$)

**likelihood**

**posterior**

$$p(\theta \mid x) = \frac{p(x \mid \theta)\, p(\theta)}{p(x)}$$

**prior**

**evidence**

- How to pick p($\theta$)?
  Pick simpler $\theta$ is better
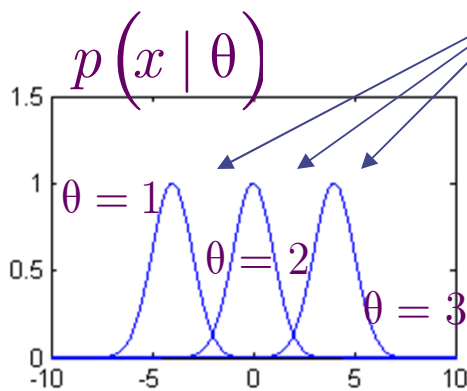  Pick form for mathematical convenience

$$p(\theta)$$



- We have data (can assume IID): $\mathcal{D} = \{x_1, x_2, \ldots, x_N\}$
- Want to get a model to compute: $p(x)$
- Want p(x) given our data... How to proceed?

# Bayesian Inference

- Want p(x) given our data... $p\left(x \mid \mathcal{D}\right) = p\left(x \mid x_1, x_2, \ldots, x_n\right)$

$$p\left(x \mid \mathcal{D}\right) = \int_\theta p\left(x, \theta \mid \mathcal{D}\right) d\theta$$

$$= \int_\theta p\left(x \mid \theta, \mathcal{D}\right) p\left(\theta \mid \mathcal{D}\right) d\theta$$ **Prior**

$$= \int_\theta p\left(x \mid \theta, \mathcal{D}\right) \frac{p\left(\mathcal{D} \mid \theta\right) p\left(\theta\right)}{p\left(\mathcal{D}\right)} d\theta$$

$$= \int_\theta p\left(x \mid \theta\right) \frac{\prod_{i=1}^{N} p\left(x_i \mid \theta\right) p\left(\theta\right)}{p\left(\mathcal{D}\right)} d\theta$$

$p\left(x \mid \theta\right)$

$\theta = 1$

$\theta = 2$

$\theta = 3$

**Many models**

**Weight on each model**

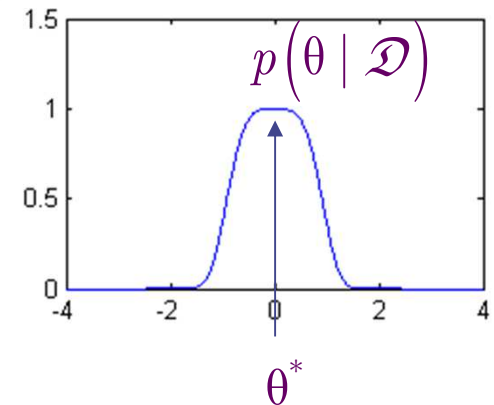$p\left(\theta \mid \mathcal{D}\right)$

# Bayesian Inference to MAP & ML

- The full Bayesian Inference integral can be mathematically tricky. MAP and ML are approximations of it…

$$p\left(x \mid \mathcal{D}\right) = \int_\theta p\left(x \mid \theta\right) \frac{\prod_{i=1}^{N} p\left(x_i \mid \theta\right) p\left(\theta\right)}{p\left(\mathcal{D}\right)} d\theta$$

$$\approx \int_\theta p\left(x \mid \theta\right) \delta\left(\theta - \theta^*\right) d\theta$$

$p(\theta \mid \mathcal{D})$

$$where \ \theta^* = \begin{cases} \arg\max_\theta \dfrac{\prod_{i=1}^{N} p\left(x_i \mid \theta\right) p\left(\theta\right)}{p\left(\mathcal{D}\right)} & MAP \\[3em] \arg\max_\theta \dfrac{\prod_{i=1}^{N} p\left(x_i \mid \theta\right) uniform\left(\theta\right)}{p\left(\mathcal{D}\right)} & ML \end{cases}$$

$p\left(\theta \mid \mathcal{D}\right)$

$\theta^*$

- Maximum A Posteriori (MAP) is like Maximum Likelihood (ML) with a prior p($\theta$) which lets us prefer some models over others

$$l_{MAP}\left(\theta\right) = l_{ML}\left(\theta\right) + \log p\left(\theta\right) = \sum_{i=1}^{N} \log p\left(x_i \mid \theta\right) + \log p\left(\theta\right)$$

# Ex: Is a coin fair?

A stranger tells you his coin is fair.

Let's assume tosses are iid with P(H)=$\mu$.

He tosses it 4 times, gets H H T H.   $\mathcal{D} = \left( H, H, T, H \right)$

What can you say about $\mu$?

# Bernoulli Probability ML

$0 \sim$ Tail

$\mu = P(H)$  $1 \sim$ Head

- Bernoulli: $p\left(x\right) = \mu^{x}\left(1-\mu\right)^{1-x}$  $\mu \in \left[0,1\right]$  $x \in \left\{0,1\right\}$

- Log-Likelihood (IID): $\sum_{i=1}^{N} \log p\left(x_i \mid \mu\right) = \sum_{i=1}^{N} \log \mu^{x_i}\left(1-\mu\right)^{1-x_i}$

- Gradient=0:

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

*N     trials/tosses*
*m    heads/1s*
*N-m  tails/0s*

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

$$\frac{\partial}{\partial \mu} \sum_{i \in class1} \log \mu + \sum_{i \in class0} \log\left(1-\mu\right) = 0$$

$$\sum_{i \in class1} \frac{1}{\mu} - \sum_{i \in class0} \frac{1}{1-\mu} = 0$$

$$m\frac{1}{\mu} - (N-m)\frac{1}{1-\mu} = 0$$

$$m\left(1-\mu\right) - (N-m)\mu = 0$$

$$m - N\mu = 0$$

$$\mu = \frac{m}{N}$$

| x=0 | x=1 |
|-----|-----|
| $\dfrac{N-m}{N}$ | $\dfrac{m}{N}$ |

# Bernoulli Bayes, Prior 1 $\quad \mathcal{D} = (H, H, T, H)$

- Assume prior μ=1/2, point mass distribution

- Posterior

$$P(\mu = r \mid \mathcal{D}) \propto P(\mathcal{D} \mid \mu = r) \times P(\mu = r)$$

- If the prior is 0 for some value, the posterior will also be 0 at that value no matter what data we see

# Bernoulli Bayes, Prior 2 $\mathcal{D} = (H, H, T, H)$

- Allow some chance of bias
- Prior $\quad$ P(μ=1/2) = 1-b

$\quad\quad\quad\quad$ P(μ=3/4) = b

- Posterior

$$P(\mu = \frac{1}{2} \mid \mathcal{D}) \propto P(\mathcal{D} \mid \mu = \frac{1}{2}) \times P(\mu = \frac{1}{2}) = \frac{1}{2^4} \times (1-b)$$

$$P(\mu = \frac{3}{4} \mid \mathcal{D}) \propto P(\mathcal{D} \mid \mu = \frac{3}{4}) \times P(\mu = \frac{3}{4}) = \left(\frac{3}{4}\right)^3 \left(\frac{1}{4}\right) \times b$$

# Bernoulli Bayes, Prior 2 $\quad \mathcal{D} = (H, H, T, H)$

- Prior $\quad$ P(μ=1/2) = 1-b

$\qquad\qquad$ P(μ=3/4) = b

- Posterior

$$P(\mu = \frac{1}{2} \mid \mathcal{D}) \propto P(\mathcal{D} \mid \mu = \frac{1}{2}) \times P(\mu = \frac{1}{2}) = \frac{1}{2^4} \times (1-b)$$

$$P(\mu = \frac{3}{4} \mid \mathcal{D}) \propto P(\mathcal{D} \mid \mu = \frac{3}{4}) \times P(\mu = \frac{3}{4}) = \left(\frac{3}{4}\right)^3 \left(\frac{1}{4}\right) \times b$$

The two are equal when $b = \dfrac{16}{43}$

# Bernoulli Bayes, Prior 3 $\mathcal{D} = (H, H, T, H)$

- Uniform prior $\mu \sim U[0,1]$

$$P(\mu = r \mid \mathcal{D}) \propto P(\mathcal{D} \mid \mu = r) \times P(\mu = r)$$

$$P(\mathcal{D} \mid \mu = r) = r^3(1-r)$$

$$\int_0^1 r^3(1-r)\,dr = \frac{1}{4} - \frac{1}{5} = \frac{1}{20}$$

*Where's the mode?*

Hence posterior $P(\mu = r \mid \mathcal{D}) = 20r^3(1-r)$

Notice for $N$ tosses with m heads, $N-m$ tails,

$$P(\mathcal{D} \mid \mu = r) = r^m(1-r)^{N-m}$$

Does this suggest a convenient prior?

# Beta Distribution

- Distribution over $\mu \in [0, 1]$

$$
\begin{aligned}
\mathrm{Beta}(\mu|a,b) &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1}(1-\mu)^{b-1} \\
\mathbb{E}[\mu] &= \frac{a}{a+b} \\
\mathrm{var}[\mu] &= \frac{ab}{(a+b)^2(a+b+1)}
\end{aligned}
$$

Here switch to typical 'sloppy' notation, using μ for variable name and its value.

# Bernoulli Bayes, Beta Prior

$$
\begin{aligned}
p(\mu|a_0, b_0, \mathcal{D}) \quad &\propto \quad p(\mathcal{D}|\mu)p(\mu|a_0, b_0) \\
&= \quad \left(\prod_{n=1}^{N} \mu^{x_n}(1-\mu)^{1-x_n}\right) \mathrm{Beta}(\mu|a_0, b_0) \\
&\propto \quad \mu^{m+a_0-1}(1-\mu)^{(N-m)+b_0-1} \\
&\propto \quad \mathrm{Beta}(\mu|a_N, b_N)
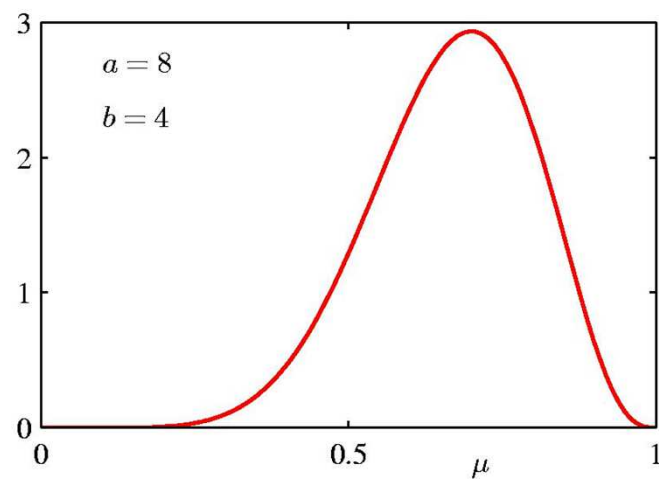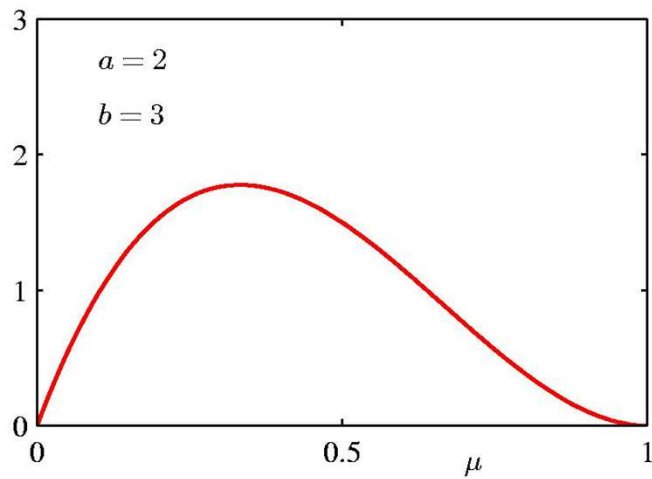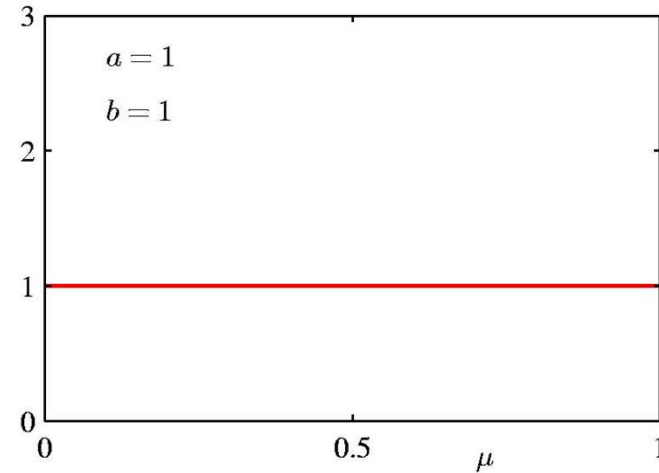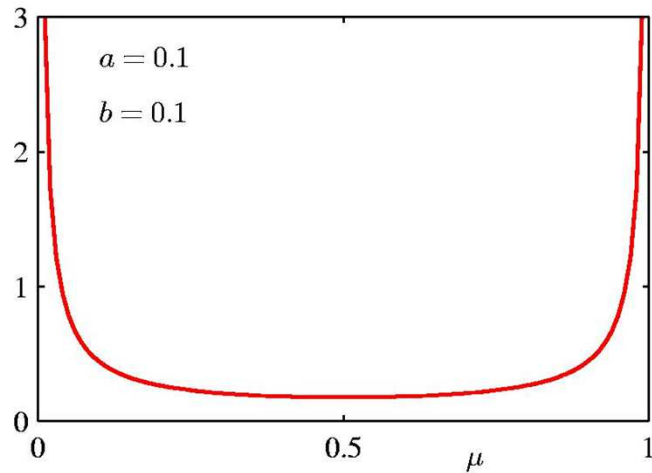\end{aligned}
$$

$$
a_N = a_0 + m \qquad b_N = b_0 + (N-m)
$$

*effective number of observations +1*

The Beta distribution provides the *conjugate prior* for the Bernoulli distribution, i.e. the posterior distribution has the same form as the prior.

All distributions in the Exponential Family (includes multinomial, Gaussian, Poisson) have convenient conjugate priors (Bishop PRML 2.4).
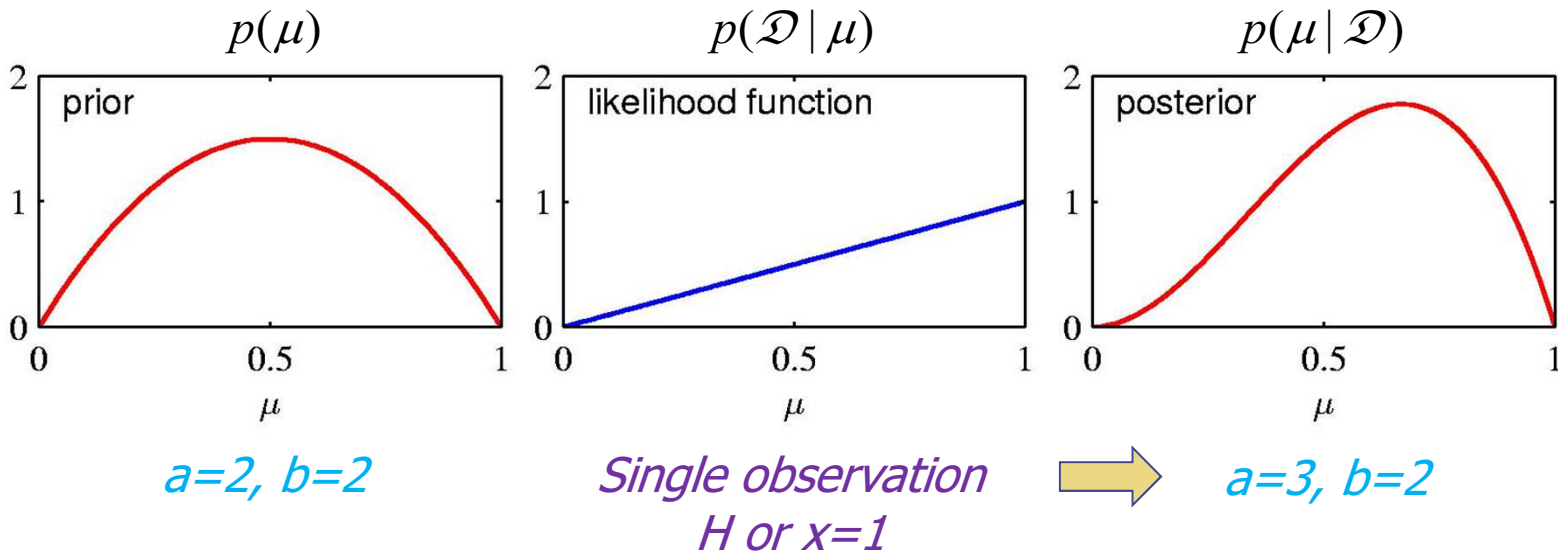
# Beta Distribution

*Which distribution is this?*

# Prior · Likelihood = Posterior

normalized

Example:



$p(\mu)$        $p(\mathcal{D}\,|\,\mu)$        $p(\mu\,|\,\mathcal{D})$

a=2, b=2        Single observation        a=3, b=2
                H or x=1

Recall our earlier example of a Uniform prior, check this works…

# Properties of the Posterior

As the size of the data set, N, grows

$$\mathbb{E}[\mu] = \frac{a_N}{a_N + b_N} = \frac{a_0 + m}{a_0 + m + b_0 + N - m} \rightarrow \frac{m}{N} = \mu_{ML}$$

$$\text{var}[\mu] = \frac{a_N b_N}{(a_N + b_N)^2 (a_N + b_N + 1)} \rightarrow 0$$

This is typical behavior for Bayesian learning.

# Prediction under the Posterior

What is the probability that the next coin toss will land heads up?

$$
\begin{aligned}
p(x = 1 | a_0, b_0, \mathcal{D}) &= \int_0^1 p(x = 1 | \mu) p(\mu | a_0, b_0, \mathcal{D}) \, \mathrm{d}\mu \\
&= \int_0^1 \mu \, p(\mu | a_0, b_0, \mathcal{D}) \, \mathrm{d}\mu \\
&= \mathbb{E}[\mu | a_0, b_0, \mathcal{D}] = \cdot \frac{a_N}{a_N + b_N}
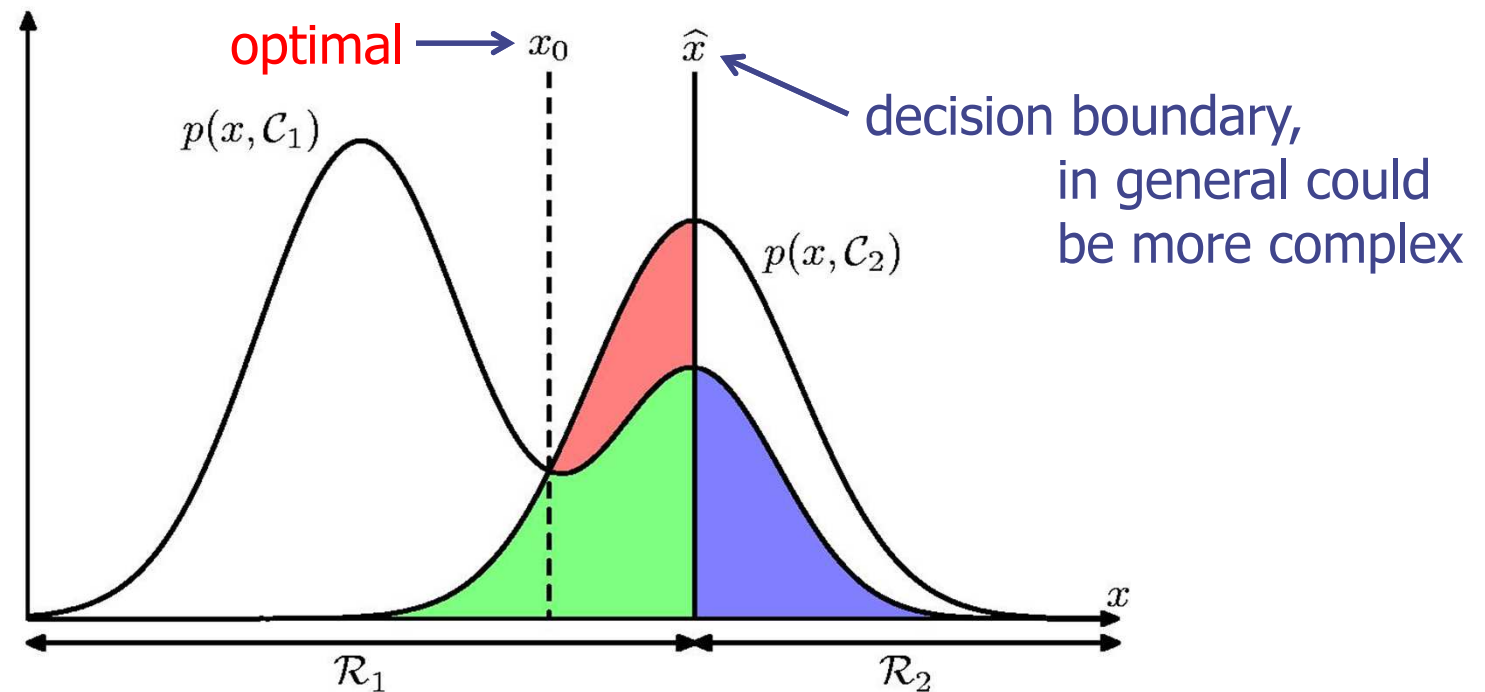\end{aligned}
$$

If we use our earlier example of a Uniform prior,

Observe $\mathcal{D} = (H, H, T, H) \Rightarrow a_N = 1 + 3, \, b_N = 1 + 1$

Now using posterior, P(next toss is a head) $= \dfrac{4}{4+2} = \dfrac{2}{3}$

# Bayesian Decision Theory

- Initially assume just 2 classes $\mathcal{C}_1, \mathcal{C}_2$
- Given input data $x$ we want to determine which class is optimal
- Various possible criteria
- Need $p(\mathcal{C}_k \mid x)$ either directly (discriminative)
- Or by Bayes, $$p(\mathcal{C}_k \mid x) = \frac{p(x, \mathcal{C}_k)}{p(x)} = \frac{p(x \mid \mathcal{C}_k) p(\mathcal{C}_k)}{p(x)}$$
  (generative)
- Divide input space into *decision regions* $\mathcal{R}_1, \mathcal{R}_2$ separated by *decision boundaries* such that $x \in \mathcal{R}_k \Rightarrow \text{assign } \mathcal{C}_k$
- How might we choose boundaries?

# Criterion 1: Minimize Misclassification Rate



optimal $\longrightarrow$ $x_0$

$\widehat{x}$

decision boundary, in general could be more complex

$p(x, \mathcal{C}_1)$

$p(x, \mathcal{C}_2)$

$\mathcal{R}_1$

$\mathcal{R}_2$

assign class 1 but should be 2    assign class 2 but should be 1

$$p(\text{mistake}) = p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1)$$

$$= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) \, d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) \, d\mathbf{x}.$$

$$p(x, \mathcal{C}_k) = p(\mathcal{C}_k \mid x) p(x)$$

# Criterion 2: Minimize Expected Loss

Example loss matrix:
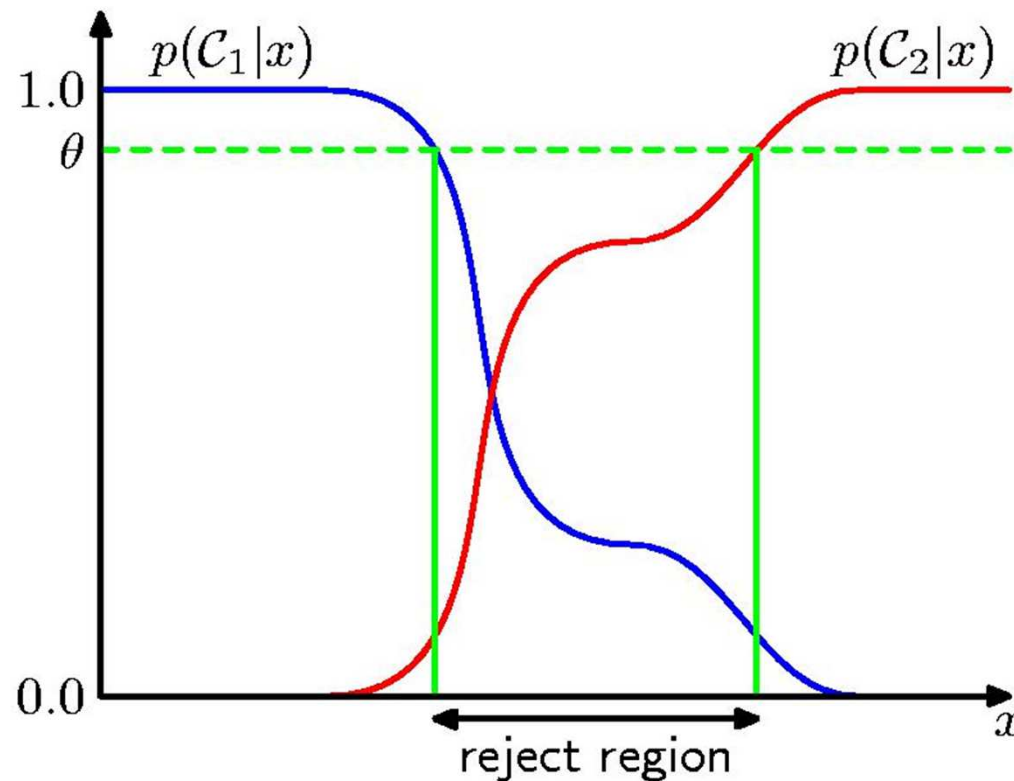
classify medical images as 'cancer' or 'normal'

$$
\begin{array}{cc}
 & \text{Decision} \\
 & \begin{array}{cc} \text{cancer} & \text{normal} \end{array} \\
\begin{array}{c} \text{Truth} \end{array}
\begin{array}{c} \text{cancer} \\ \text{normal} \end{array}
& \left( \begin{array}{cc} 0 & 1000 \\ 1 & 0 \end{array} \right)
\end{array}
$$

Now $\quad \mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) \, \mathrm{d}\mathbf{x}$

Choose regions $\mathcal{R}_j$ to minimize Expected Loss

$$\mathbb{E}[L] = \sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

# Reject Option



If $p(C_1 \mid x) \approx p(C_2 \mid x)$ then less confident about assigning class.

One possibility is to reject or refuse to assign.