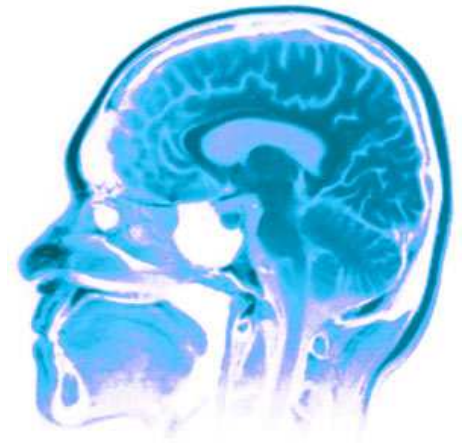# CPSC340

# Bayesian learning

Nando de Freitas

*October, 2012*

*University of British Columbia*
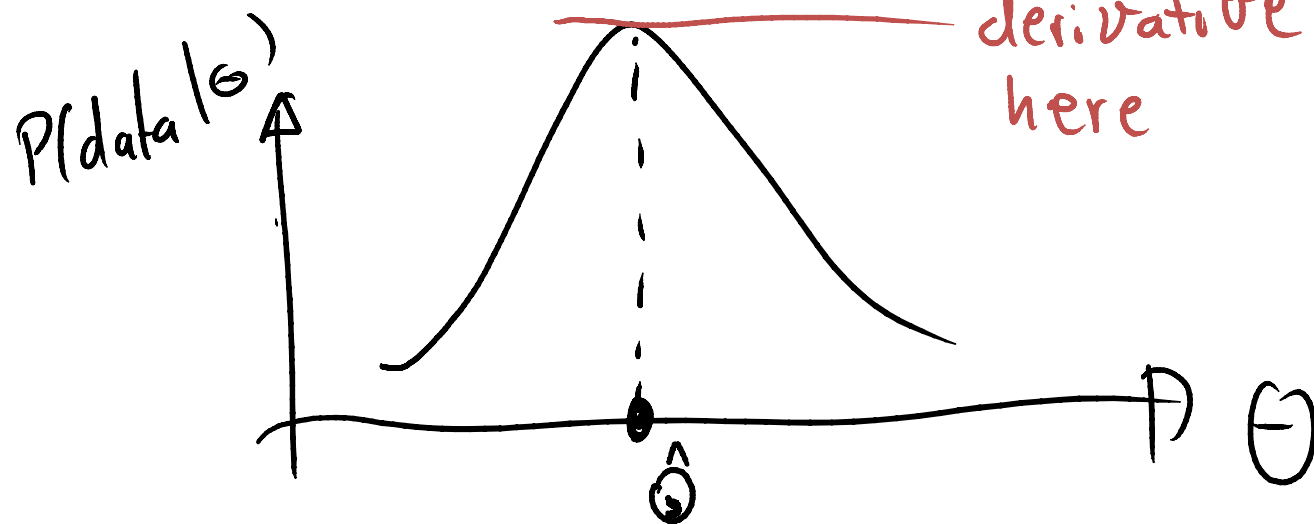
# Maximum likelihood revision

Find $\Theta$ by maximizing $P(\text{data} \mid \Theta)$.

e.g. for a coin $P(\underbrace{\text{data}}_{X_{1:n}} \mid \Theta) = \Theta^m (1-\Theta)^{n-m}$
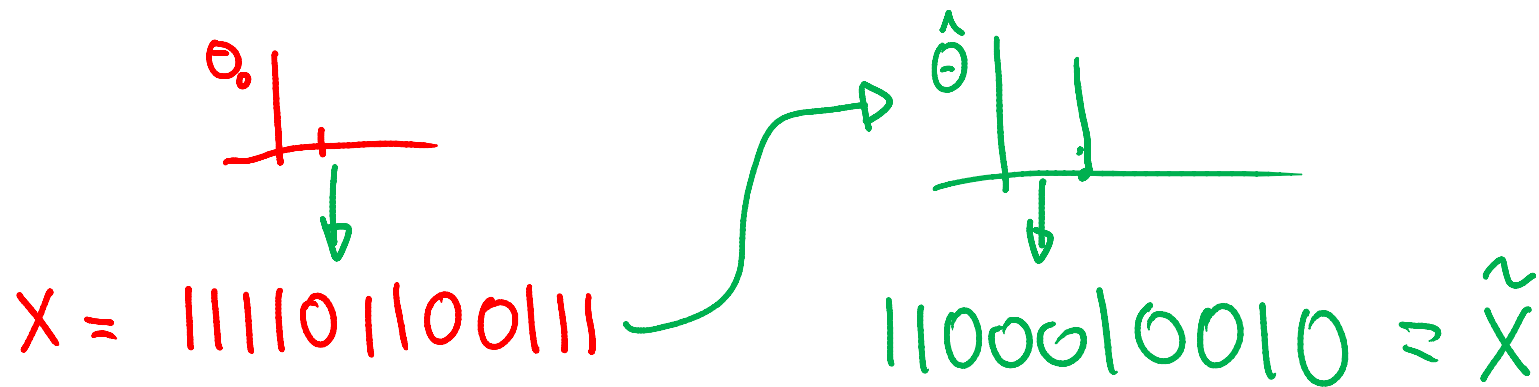
$m = \#\ 1\text{'s}$

$n = \#\ 1\text{'s and } 0\text{'s}$

Slope derivative is zero here

# Maximum likelihood revision

$\hat{\Theta}_{ML}$ is the theta that minimizes

$$error(\theta) = Information(\Theta_0) - Information(\theta)$$

true theta

any theta



$$X = 11111011001111$$

$$11000100010 = \tilde{X}$$

choose $\hat{\Theta}$ to minimize $|f(x) - f(\tilde{x})|$

# Outline of the lecture

This lecture introduces us to our second strategy for learning: **Bayesian learning**. The goal is for you to learn:

- ❑ Definition Beta prior.
- ❑ How to use Bayes rule to go from **prior beliefs** and the **likelihood of the data** to **posterior beliefs**.

# Bayesian learning procedure

**Step 1:** *Given **n** data, $x_{1:n} = \{x_1, x_2, ..., x_n\}$, write down the expression for the likelihood:*

$$p(x_{1:n} \mid \theta) = \theta^m (1-\theta)^{n-m} \quad \text{(for a coin)}$$

**Step 2:** *Specify a prior: $p(\theta)$*

$$P(\theta \mid x_{1:n}) = \frac{1}{\text{Const}} P(x_{1:n} \mid \theta) P(\theta)$$

**Step 3:** *Compute the posterior:*

$$p(\theta \mid x_{1:n}) = \frac{p(x_{1:n} \mid \theta) \, p(\theta)}{p(x_{1:n})}$$

Notation:

$$p(\theta \mid x_{1:n}) \propto p(x_{1:n} \mid \theta) \, p(\theta) \quad \leftarrow$$

↳ proportional to

# Bayesian learning procedure

$$P(A) = \sum_B P(AB)$$

$$P(A) = \int P(AB)\, dB$$

***Posterior:*** *Compute the posterior:*

$$p(\theta \mid x_{1:n}) \propto p(x_{1:n} \mid \theta)\, p(\theta)$$

$$\int P(\theta \mid x_{1:n})\, d\theta = 1 \qquad P(\theta \mid x_{1:n}) = \frac{P(x_{1:n} \mid \theta)\, P(\theta)}{\int P(x_{1:n} \mid \theta)\, P(\theta)\, d\theta}$$

***Marginal likelihood:*** $\quad p(x_{1:n}) = \int P(x_{1:n} \mid \theta)\, P(\theta)\, d\theta$

# Bayesian learning for coin model

***Step 1:*** *Write down the likelihood of the data (i.i.d. Bernoulli in our case):*

$$p(x_i / \theta) = \theta^{x_i} (1-\theta)^{1-x_i}$$

$$p(x_{1:n} / \theta) = \theta^{m} (1-\theta)^{n-m}$$

$$x_i \in \{0, 1\}$$

$$m = \# \ 1\text{'s}$$

# Bayesian learning for coin model

**Step 2:** *Specify a prior on* **θ**. *For this, we need to introduce the Beta distribution.*

We know $\theta$ is continuous and $0 \leq \theta \leq 1$

$\theta = P(X_i = 1) \quad \forall i$

$P(\theta)$?

$$P(\theta) = \left[ \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\,\Gamma(\beta)} \right] \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

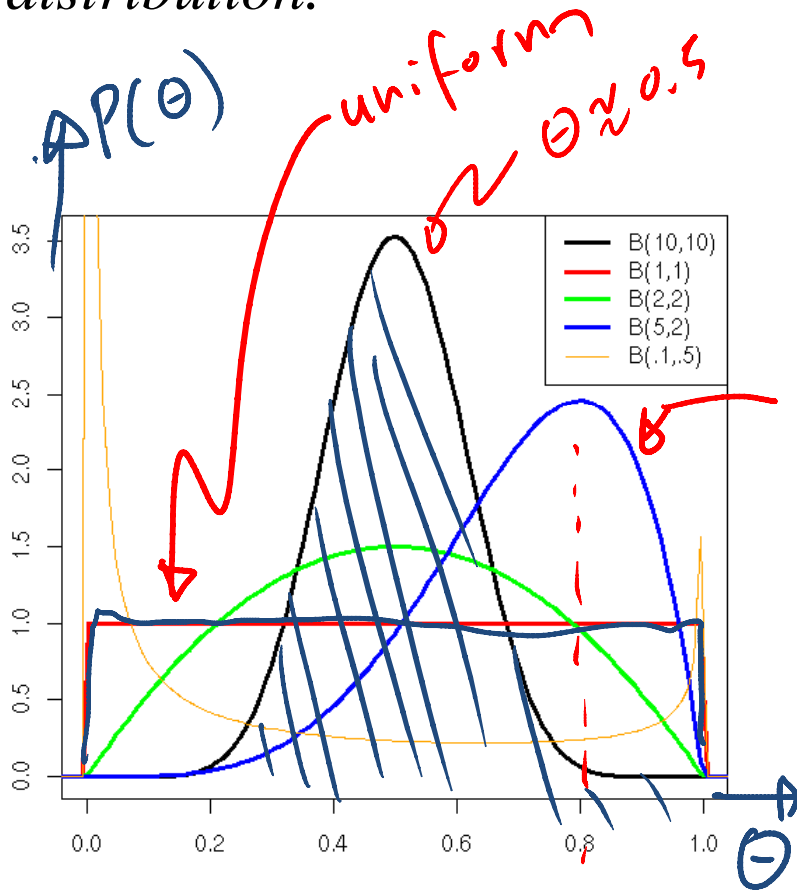$(\alpha, \beta)$ hyperparams

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

$$\int \theta^{\alpha-1}(1-\theta)^{\beta-1}\, d\theta = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

$$\log(z) = \int_1^z \frac{1}{x}\, dx$$

$$\int P(\theta)\, d\theta = 1 = \int \frac{\Gamma(\alpha+\beta)\,\theta^{\alpha-1}(1-\theta)^{\beta-1}}{\Gamma(\alpha)\Gamma(\beta)}\, d\theta$$

# Bayesian learning for coin model

*Step 2: Specify a prior on θ. For this, we need to introduce the Beta distribution.*



uniform

$\theta \approx 0.5$

$P(\theta)$

Legend:
- B(10,10)
- B(1,1)
- B(2,2)
- B(5,2)
- B(.1,.5)

$\theta$

0.8

If $\alpha = 1$ $\beta = 1$

$$P(\theta) \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$$

$$= \theta^0(1-\theta)^0$$

$$= 1$$

ie. $P(\theta) \propto 1$ (uniform)

$$\text{mean}(\theta) = \frac{\alpha}{\alpha+\beta}$$

# Bayesian learning for coin model

*Step 3: Compute the posterior:*

$$p(\theta \mid x_{1:n}) \propto p(x_{1:n} \mid \theta)\, p(\theta) =$$

$$= \theta^m (1-\theta)^{n-m}\ \theta^{\alpha-1}(1-\theta)^{\beta-1}$$

$$= \theta^{m+\alpha-1}(1-\theta)^{n-m+\beta-1}$$

$$= \theta^{\alpha'-1}(1-\theta)^{\beta'-1} \qquad \alpha' = m+\alpha$$
$$\beta' = n-m+\beta$$

$$p(\theta \mid x_{1:n}) = \frac{1}{const}\ \theta^{\alpha'-1}(1-\theta)^{\beta'-1}$$

$$= \frac{\Gamma(\alpha'+\beta')}{\Gamma(\alpha')\Gamma(\beta')}\ \theta^{\alpha'-1}(1-\theta)^{\beta'-1}$$

# Example

*Suppose we observe the data, $x_{1:6} = \{1, 1, 1, 1, 1, 1\}$, where each $x_i$ comes from the same Bernoulli distribution (i.e. it is independent and identically distributed (iid)). What is a good guess of $\theta$?*

$$\hat{\Theta}_{ML} = 1$$

*We can compute the posterior and use its mean as the estimate.*

$$P(\theta | X_{1:6}) \propto \theta^{6+\alpha-1} (1-\theta)^{0+3-1}$$

1000000 1's

$$\hat{\Theta}_B \approx 1$$

$$\hat{\Theta}_B = \mathbb{E}(\theta | X_{1:6}) = \frac{6+\alpha}{6+\alpha+3}$$

*Using a prior Beta(2,2):*

$$\hat{\Theta}_B \approx \frac{8}{10} = 0.8$$

$$Beta(1,1)$$

$$\hat{\Theta}_B = \frac{7}{7+1} = \frac{7}{8}$$

*Using a prior Beta(1,0.01):*

$$\hat{\Theta}_B \approx \frac{7}{7+0.01} \approx 1$$

# Next lecture

In the next lecture, we apply our learning strategies to Bayesian networks.