# Bayesian Model Selection

Last time: bias & variance as a way to understand how model complexity affects generalization.

(freq: data <u>are</u> random    $\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)]$

Today: <u>Bayesian view</u>: data are fixed; params <u>are</u> random
$\downarrow$

is a formula for many things

1) compute $P(w|X,Y)$ } <u>posterior</u> over params

2) compute $P(y^* | x^*, X, Y) = \int_w P(y^*|x^*, w) \, p(w|X,Y)$

<u>posterior predictive</u> { ↗ prediction   ↑ some new point

(BAYES is all about model-averaging)

3) compute $P(Y|X) = \int_w P(Y|X,w) \, p(w)$ } <u>marginal likelihood</u>

## Posterior over params:

Bayes rule: $P(w|X,Y) = \dfrac{P(Y|X,w) \, P(w|x)}{P(Y|X)}$

← assume prior over w doesn't dep. on X

← notice: this is a constant wrt w.

$\propto P(Y|X,w) \, P(w)$

$\propto P(w,Y|X)$

depending on $P(Y|X,w)$, $P(w) \to$ easy or very hard
$\downarrow$
conjugate.

## Example w/ Beta-Bernoulli Model

Model: coin come up "1" with prob. $\theta$ } x is the flip

$p(x|\theta) = \theta^x (1-\theta)^{1-x}$    likelihood

$p(x_1 \cdots x_N | \theta) = \prod_{n=1}^{N} \theta^{x_n} (1-\theta)^{1-x_n}$

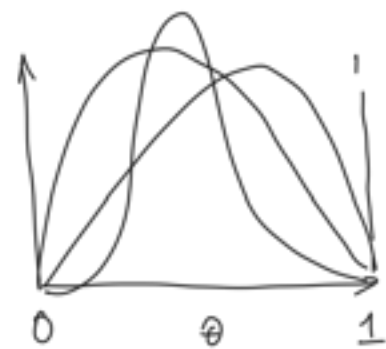$\qquad = \theta^{n_1} (1-\theta)^{n_0}$ , where $n_1 + n_0 = N$

Aside: from this, we could compute $\theta_{MLE}$, $\theta$ maximum likelihood,

$\max_\theta \log p(x_1 - x_n | \theta) = \max_\theta n_1 \log \theta + n_0 \log(1-\theta)$

$\Rightarrow \theta_{MLE} = \dfrac{n_1}{n_0 + n_1}$

Now, let's put a prior on $\theta$. Beta distribution: $p(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \; \mathbb{I}(0,1)$

for $\alpha, \beta > 0$

$(\theta$ in $[0,1])$

for $\alpha, \beta > 1$
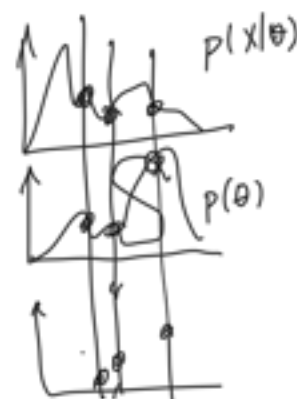
if $\alpha, \beta = 1 \Rightarrow$ uniform

if $\alpha, \beta < 1$
"sparsity favoring"

Why this choice?

$$p(\theta | X) \propto p(X|\theta)\, p(\theta)$$
$$\propto \frac{\theta^{n_1}(1-\theta)^{n_0} \; \theta^{\alpha-1}(1-\theta)^{\beta-1}}{}$$
$$= \theta^{n_1+\alpha-1}(1-\theta)^{n_0+\beta-1} \quad \text{looks like another Beta!!}$$

$\text{Beta}(\alpha, \beta) \longrightarrow \boxed{\text{Beta}(n_1+\alpha,\; n_0+\beta)}$
$\quad p(\theta) \qquad\qquad\qquad p(\theta|X)$

$p(X|\theta)$

$p(\theta)$

$\theta$

fill in some data:
- $\alpha = 1$, $\beta = 1$
- we see 2 heads: $n_0 = 0$, $n_1 = 2$
- new Beta$(3, 1)$

With this dist, we can compute M.A.P. (last time)

$$\max_{\theta} \log p(X|\theta)\, p(\theta) \Rightarrow \boxed{\frac{n_1 + \alpha - 1}{n_1 + n_0 + \alpha + \beta - 2}} = \theta_{MAP} \quad \text{mode of Beta}$$

with this dist, we can also compute posterior predictive

$$p(x=1 \mid x_1 \cdots x_N) = \int_\theta \underbrace{p(x=1|\theta)}_{\text{predict}} \underbrace{p(\theta|x_1\cdots x_N)}_{\text{posterior}}$$

$$= \int \theta \, p(\theta|x_1\cdots x_N) = E_{p(\theta|x_1\cdots x_N)}[\theta] = \frac{\alpha + n_1}{\alpha + \beta + n_1 + n_0}$$

mean of Beta

$\theta = \frac{n_1}{\quad} = 1 \quad$ argmax $p_\theta(x|\theta)$

$$\theta_{MLE} = \frac{n_1}{n_1 + n_0} = 1 \qquad \text{argmax} \quad$$

$$\theta_{MAP} = 1 \qquad \text{argmax} \quad P(X|\theta) P(\theta)$$

$$\theta_{PP} = 3/4 \qquad \underset{\theta}{\mathbb{E}} \; E_{p(\theta|x_1 \cdots x_N)}[\theta]$$

in a diff problem: we could compare

$$P(y|x, \theta_{MLE})$$
$$P(y|x, \theta_{MAP})$$
$$P(y|x, X, Y) \leftarrow PP$$

finally: $\underbrace{P(x_1 \cdots x_N)}_{P(X)} = \underbrace{\int \underbrace{\theta^{n_1}(1-\theta)^{n_0}}_{P(X|\theta)} \underbrace{\theta^{\alpha-1}(1-\theta)^{\beta-1} \boxed{Z_{\alpha,\beta}}}_{P(\theta)}}$

looks like
Beta$(\alpha+n_1, \beta+n_0)$

$\dfrac{Z_{\alpha+n_1, \beta+n_0} \cdot}{Z_{\alpha+n_1, \beta+n_0} \cdot}$

$$= \underbrace{\frac{Z_{\alpha,\beta}}{Z_{\alpha+n_1, \beta+n_0}}}_{\text{marginal likelihood}} \underbrace{\int \text{Beta}(\alpha+n_1, \beta+n_0)}_{1}$$

$P(\theta)$
Beta$(\alpha, \beta)$
vs
$P(\theta)$

Why is this marginal likelihood useful?

Bayesian Occam's Razor : $P(X)$ — implicitly $P(X| \underset{\text{class}}{\text{model}})$



$p(\text{model})$ — $p(\text{model for "simple" class})$

model needed to explain data

$p(\text{model})$

linear   poly k   poly k+1

models : less to more complex $\longrightarrow$
(linear vs. poly$_k$ vs. poly w/ bigger k)

1

---

Concept Check



(1,1)

(0,0)

$y = f(x) + \varepsilon$
$\underbrace{\quad}_{\substack{\text{assume} \\ \text{tiny}}}$

1. What parabolas can we fit to these data, assume noise is very small.
   (what is implication of small noise)

   $\rightarrow$ parabola will go through data

2. $y = a_0 + a_1 x + a_2 x^2$ can we define a single parabola that fits data?
   $a_0 = 0$
   $a_2 = 1 - a_1$

$1 = a_1 + a_2$

3. improper prior on $a_2$ $[\ p(a_2) = 1 \ \text{for all } a_2\ ]$ , and unif $[-1, 1]$ for $a_1$.

What is posterior predictive at $x = 1/2$.

if $a_1 = 1$: $a_2 = 0 \to 0$
if $a_1 = -1$: $a_2 = 2 \to 1/2$

$$y(x=\tfrac{1}{2}) = \tfrac{1}{4} + \tfrac{1}{4}a_1 \quad \Rightarrow \text{avg } \tfrac{1}{4}$$

4. What would be different if we assumed a linear model $y = a_0 + a_1 x$?
   (marginal likelihood)

   linear model: only one model

   posterior pred: $1/2$

   marg lik: