# external validation of SBI (against real data, contd.)
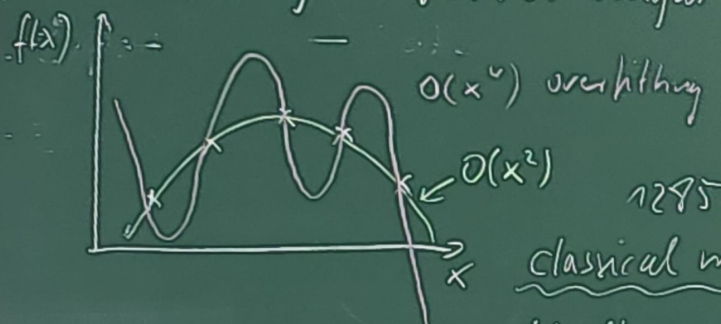
- **model misspecification detection**: is $X^{oss}$ an outlier to the simulation?

  if yes $\rightarrow$ reject $X^{oss}$ and answer "I don't know"

  a. exploit the feature detection network:
    - add loss $\propto$ MMD$(p(h(x)) \mid N(0,\mathbb{I}))$ to pull summary/feature distr. towards standard normal
    - reject $X^{oss}$ if $h(X^{oss})$ is an outlier of $N(0,\mathbb{I})$

- **model comparison & selection** (between competing theories)

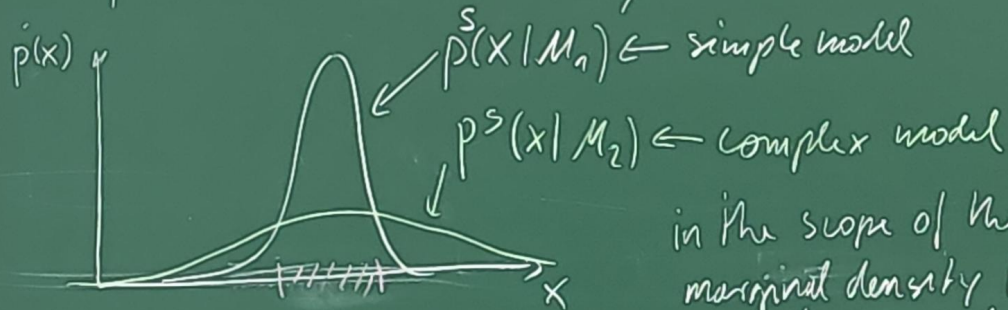  - measuring training error might not be enough, because overfitting might occur



$f(x)$

$O(x^4)$ overfitting

$O(x^2)$

- need trade-off between model accuracy & complexity

(intuition: simpler models generalize better)

1285-1347: "Occam's razor" prefer simplest theory

classical model selection criteria:

Akaike criterion $AIC = 2\,E[NLL] + 2\,size$ $\leftarrow$ # model paramet. (e.g. # dimensions of linear model)

Bayesian information crit    $BIC = 2\,\mathbb{E}[NLL] + \log(N)$  size

- in Bayesian inverence, model complexity - is            dataset size
  penalized automatically

$p(x)$

$\overset{S}{p}(X|M_1) \leftarrow$ simple model        $M_1, \ldots, M_L$  competing
                                                              simulations

$p^S(x|M_2) \leftarrow$ complex model



in the scope of the simple model ⊬⊬⊬⊬I, the
marginal density for $M_1$ is much higher than for $M_2$
due to normalization of probs

$\Rightarrow$ it $X^{obs} \in$ |⊬⊬⊬|, during training of SBI, it typically came from $M_1 \Rightarrow$ automatically
                prefer $M_1$ during inference as well                        noise outsourcing

forward model :    $\underbrace{p(M)}_{\sim uniform(1,L)} \cdot p(Y|M) \cdot p(\eta|M,Y) \cdot \delta\!\left(X - \phi_M(Y,\overset{\kappa}{\eta})\right)$

                                                                              simulation M
                   $\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad}$
                            $= p(X|Y,M)$    likelihood

marginal density $\qquad p(X|M) = \int p(X|Y,M)\, p(Y|\mu)\, dy$

$\Rightarrow$ model comparison ① : Bayes factor $= \dfrac{p(X = x^{oss} | M_\ell)}{p(X = x^{oss} | M_\ell')} \begin{cases} > 1 & \text{prefer } M_\ell \\[2mm] < 1 & \text{prefer } M_{\ell'} \end{cases}$

posterior for model preferences

$$p(M|X) = \frac{p(X|M)\, p(M)}{p(X)} \qquad p(x) = \sum_{\ell=1}^{L} p(M = \ell)\, p(x | M = \ell)$$

$\Rightarrow$ model comparison ② : posterior odds $\dfrac{p(M_\ell | X = x^{obs})}{p(M_{\ell'} | X = x^{obs})}$ $\qquad$ equal to Bayes factor if $p(M) = \text{uniform}(1, L)$

$\cdot$ practical alg. : train a standard softmax classifier for $p(M|X)$

$\bullet$ comparison thresholds $\dfrac{p(M_\ell | X^{oss})}{p(M_{\ell'} | X^{obs})} = r \begin{cases} \frac{1}{3} < r < 3 & : \text{no significant differences} \\[1mm] 3 < r < 10 & : \text{substantial evidence for } M_\ell \\[1mm] 10 < r < \frac{30}{100} & \quad\quad \text{strong} \quad -//- \\[1mm] \frac{30}{100} < r & \quad\quad \text{overwhelming} \; -//- \\[1mm] \text{likewise} \; \frac{1}{30}, \frac{1}{10}, \frac{1}{3} & : \text{evidence for } M_{\ell'} \end{cases}$

## external validation alg.

given: competing theories $M_1, \ldots, M_L$

[ epidemiology: different # compartments, priors, observation uncertainly etc ]

① create synthetic training data

$$TS_\ell = \left\{ \left( Y_{\ell i} \sim p(Y|M=\ell), \ X_n \sim p(X|Y_{n\ell}, M=\ell) \right) \right\}_{i=1}^{N_\ell} \qquad TS = \{ TS_1, \ldots, TS_L \}$$

② train a separate SBI model for each $TS_\ell$ with MMD, so that $p(h_\ell(x)) = N(0, \mathbb{I})$

③ perform internal validation for each $\ell$, redisign $SBI_\ell$ until successful

④ train a softmax classifier $p(M|X)$ using combined $TS$

cross-entropy
loss

$$\hat{p}(M|X) = \arg\min_{p} \frac{1}{L} \sum_{\ell=1}^{L} \frac{1}{N_\ell} \sum_{n=1}^{N_\ell} - \log p(M=\ell|X=X_{n\ell})$$

⑤ external validation: given $X^{obs}$

ⓐ model misspecification detection: $M^{in} = \{ \ell : h_\ell(X^{obs}) \text{ is inlier of } N(0, \mathbb{I}) \}$

ⓑ compute logits of model classifier $S_\ell$ ( penultimate layer, before softmax )

ⓒ define classifier: $p(M|X^{obs}) = softmax(S_\ell : \ell \in M^{in})$
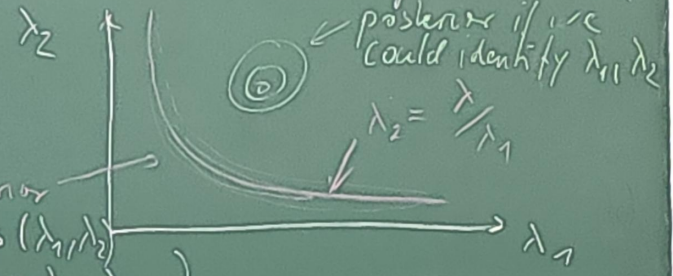
ⓓ model comparison by posterior odds

# parameter degeneracy

- sometimes, some elements in $Y$ cannot be fully identified from $X$
  $\Rightarrow$ correlations in the posteriors

- example — epidemiologist write SIR equations in terms of natural/conceptional parameters
  - $\lambda_1$ average number of people a healthy person meets per day
  - $\lambda_2$ fraction of dangerous meetings leading to transmission
  - in SIR eq., we always have $\lambda_1 \cdot \lambda_2 \Rightarrow$ we cannot distinguish them
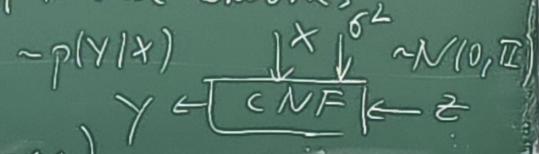  - but, we can infer $\lambda = \lambda_1 \cdot \lambda_2$

$\Rightarrow$ if we still use $\lambda_1$ & $\lambda_2$ $\Rightarrow$ posterior shows the dependency

actual posterior
(infinitely many pairs $(\lambda_1, \lambda_2)$
for fixed $\lambda = \lambda_1 \cdot \lambda_2$ )

posterior if we could identify $\lambda_1, \lambda_2$

$\lambda_2 = \lambda / \lambda_1$

- here, cause of degeneracy is easy to spot, but generally difficult
  and hard to distinguish from bad convergence of neural networks
- CNF struggle when $p(Y|X)$ is degenerate, because $\sim p(Y|X)$
  - code distribution $N(0, \mathbb{I})$ has $D$ dimensions
  - but $p(Y|X)$ has $< D$ dimensions ($\hat{=}$ degeneracy)

$X \downarrow \sigma^2 \sim N(0, \mathbb{I})$
$Y \leftarrow \boxed{CNF} \leftarrow Z$

theorem: <u>bijective transforms</u> are only possible if dimension does not change

e.g. NF

- trick to learn a good approximation    <u>Soft Flow</u>  [ Kim et al. 2020 ]

  - idea: add <u>noise</u> to $Y_i$ from TS to make it D-dimensional

    D-dimension Gaussian $N(0, \sigma^2 \mathbb{I})$

  - vary $\sigma^2$ during training according to $\sigma^2 \sim p(\sigma^2)$

  - tell NF about current value of $\sigma^2$          $\overbrace{\phantom{xxxxx}}^{\text{some prior}}$
    ( additional condition    $p(Y | X, \sigma^2)$ )

    $\Rightarrow$ network learns to generate data with given amount of noise $\sigma^2$

  - at inference time, make $\sigma^2 \rightarrow \sigma^2_{min}$ ( $\sigma^2_{min}$ smallest prior value during training )

  $\Rightarrow$ line $\lambda_2 = \lambda / \lambda_1$ becomes as narrow as possible

  ( ideally, do inference $\sigma^2 \rightarrow 0$, but practical NF saturate at some finit $\sigma^2_{min}$ )

  [ do not confuse with Noise Net  adds noise to $X$
              Soft Flow    — || —    $Y$ ]

$Y_2$

cannot be learned by NF

can be learned

$Y_1$