# Model Comparison

- PPCs for checking model fit. But how do we compare different models (regardless of fit)?

- Compare models by their out-of-sample predictive power

- Leave-one-out cross validation

  - Fit model $n$ times, leaving out one observation

  - Evalute predcitive accuracy on left out observations

  - Typically used for parameter tuning

- Can use ideas from LOO-CV to compare models

# Measures of Predictive Accuracy

- Point prediction: predict a single unknown future observation

- Measures of predictive accuracy are called *scoring functions* (e.g. MSE)

- Probabilistic predictions account for uncertainty

  - **Scoring rules**: measure of accuracy based on probabilistic predictions

  - Examples: quadratic, logarithmic and zero-one scores

# Log-likelihood as a scoring rule

- As a measure of accuracy, we call it the "log predictive density"

- $\log(p(y|\theta))$ is proportional to MSE if the data are normal with constant variance

  *chosen Model*

- Connection to Kullback-Leibler divergence

  - KL: $-E_p log(q(y)/p(y|\theta)) = -E_p log(q(y)) + E_p log(p(y))$

    *Fwd-KL*

    *True Model (unknown)*

  - Asymptotically, the lowest KL model is the one with the highest expected log predictive density

# Out-of-sample prediction

Ideal: Predict fit on new data.

$$\log P_{post}(\tilde{y}) = \log E_{post} P(\tilde{y}|\theta) = \int P(\tilde{y}|\theta) P(\theta|y) \, d\theta$$

# Expected Log Pointwise Predictive Density

$(ELPPD)$ $(one\ obs)$

$$E_f \log P_{post}(\tilde{y}) = \int \log P_{post}(\tilde{y}) f(\tilde{y}) d\tilde{y}$$

True Model $f(\tilde{y})$

$$ELPD: \sum E_f \log P_{post}(\tilde{y})$$

$$\widehat{ELPD}_{LOO} : \sum_{i=1}^{n} \log P(y_i | y_{-i})$$

$$P(y_i | y_{-i}) = \int P(y_i | \theta) P(\theta | y_{-i}) d\theta$$

problem: Computation

Need: $P(\theta | y_{-i}) \quad \forall_i$

Have: $P(\theta | y)$

Importance Sampling

$$P(y_i | y_{-i}) = \int P(y_i | \theta) P(\theta | y_{-i}) d\theta =$$

$$\int P(y_i | \theta) \overbrace{\left( \frac{P(\theta | y_{-i})}{P(\theta | y)} \right)}^{\text{I.W.}} P(\theta | y) \, d\theta$$

When $y_{1,\dots,} y_n$ are iid.

$$\left( \frac{P(\theta | y_{-i})}{P(\theta | y)} \right) = \frac{1}{P(y_i | \theta)}$$

$$W_{is} = \frac{1}{P(y_i | \theta^s)}$$

# LOO Importance Sampling

- When the data are conditionally independent,

$$r_i^s = \frac{1}{p(y_i|\theta^s)} \propto \frac{p(\theta^s|y_{-i})}{p(\theta^s|y)} \text{ to get importance sampling}$$

estimates

- $p\left(\tilde{y}_i \mid y_{-i}\right) \approx \dfrac{\sum_{s=1}^{S} r_i^s p\left(\tilde{y}_i|\theta^s\right)}{\sum_{s=1}^{S} r_i^s}$

- $p\left(y_i \mid y_{-i}\right) \approx \dfrac{1}{\frac{1}{S} \sum_{s=1}^{S} \frac{1}{p(y_i|\theta^s)}}$

https://link.springer.com/article/10.1007/s11222-016-9696-4

# Pareto Smoothed Importance Sampling

1. Fit the generalized Pareto distribution to the 20% largest importance ratios $r_s$

2. Replacing the 20% largest ratios by the expected values of the order statistics of the fitted generalized Pareto distribution

3. Truncate the weights to ensure finite variance (see paper)

The above steps must be performed for each data point i.

$$\widehat{\text{elpd}}_{\text{psis}-\text{loo}} = \sum_{i=1}^{n} \log\left( \frac{\sum_{s=1}^{S} w_i^s p\left(y_i | \theta^s\right)}{\sum_{s=1}^{S} w_i^s} \right)$$

```
1  library(loo)
2  nb_model <- cmdstan_model("nb_model.stan")
3  nb_fit = nb_model$sample(
4                  data=list(n1=num_bachelors, n2=num_no_bachelors,
5                  y1=bachelors_data, y2=no_bachelors_data),
6                  refresh=0)
```

Running MCMC with 4 parallel chains...

Chain 1 finished in 0.4 seconds.
Chain 2 finished in 0.4 seconds.
Chain 3 finished in 0.4 seconds.
Chain 4 finished in 0.4 seconds.

All 4 chains finished successfully.
Mean chain execution time: 0.4 seconds.
Total execution time: 0.5 seconds.

```
1  loo_compare(list("pois"=pois_fit$loo(),
2                   "zip"=zip_fit$loo(),
3                   "nb"=nb_fit$loo()))
```

```
      elpd_diff se_diff
zip    0.0       0.0
nb    -3.2       2.4
pois  -8.6       4.6
```

# Galaxies Example

- Velocities in km/sec of 82 galaxies (Corona Borealis region).

- Multimodality in such surveys is evidence for voids and superclusters in the far universe.

- Statistical question: how many clusters are there in this dataset?

# Model Comparison

```r
1  library("loo")
2  galaxy_speeds <- log(MASS::galaxies)
3  galaxy_results <- list()
4
5  ## Run model for each of k clusters
6  mix_model <- cmdstanr::cmdstan_model("mix_model.stan")
7  for(k in 1:5) {
8    galaxy_results[[k]] <-
9      mix_model$sample(
10       data=list(K=k, N=length(galaxy_speeds), y=galaxy_speeds),
11       refresh=0, show_messages=FALSE)
12 }
```

Warning: 63 of 4000 (2.0%) transitions ended with a divergence.
See https://mc-stan.org/misc/warnings for details.

Warning: 16 of 4000 (0.0%) transitions hit the maximum treedepth limit of 10.
See https://mc-stan.org/misc/warnings for details.

# Model Comparison

```r
loo_compare(galaxy_results[[1]]$loo(),
            galaxy_results[[2]]$loo(),
            galaxy_results[[3]]$loo(),
            galaxy_results[[4]]$loo(),
            galaxy_results[[5]]$loo())
```

```
       elpd_diff se_diff
model3    0.0       0.0
model5   -0.4       1.4
model4   -1.0       0.7
model2  -10.7       3.7
model1  -41.1       8.7
```

# Model Comparison

```r
1  galaxy_results[[4]]$loo()
```

```
Computed from 4000 by 82 log-likelihood matrix

          Estimate    SE
elpd_loo      32.0   9.5
p_loo          9.2   1.3
looic        -64.0  18.9
------
Monte Carlo SE of elpd_loo is 0.1.

Pareto k diagnostic values:
                          Count Pct.    Min. n_eff
(-Inf, 0.5]    (good)       81   98.8%   160
 (0.5, 0.7]    (ok)          1    1.2%   953
  (0.7, 1]     (bad)         0    0.0%   <NA>
```