

Bayesian multimodeling: Bayesian inference

MIPT

2023

Coin problem

A person flips a coin N times. What's the probability of getting tails on a coin?

Coin problem

A person flips a coin 3 times. All 3 times it comes up tails. What's the probability of getting tails on a coin?

Naive approach

$$\mathbf{X} = [1, 1, 1];$$

$$x \sim \text{Bin}(w);$$

$$\hat{w} = \arg \max_p L(\mathbf{X}, w);$$

$$\rightarrow \hat{w} = 1.$$

Challenge: three events are not enough to estimate the distribution of heads and tails.

Frequentist and Bayesian statistics

Frequentist statistics

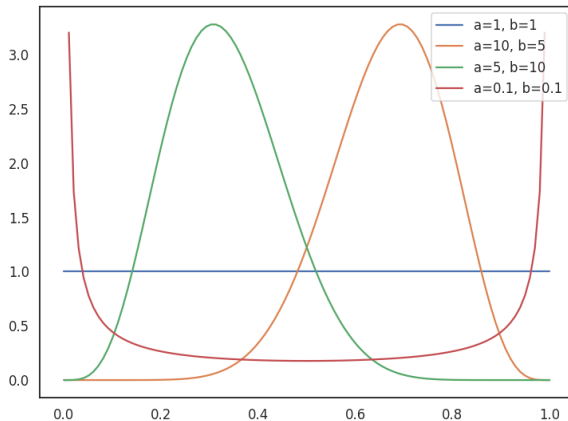
- Model parameter is a constant that is required to be estimated
- Probability is estimated purely from event frequency

Bayesian statistics

- Model parameter is a random value
 - ▶ We cannot “estimate” random value
 - ▶ But we can estimate its distribution parameters
- Probability is estimated wrt our prior beliefs about data and parameter distribution
 - ▶ The more data we get the closer our estimation to MLE
 - ▶ In general our estimation is strongly relies on the prior

Beta-distribution: recap

- corresponds to the *prior* beliefs about Bernoulli distribution
- interpretation: “effective number of events $w = 1, w = 0$ ”
- With $n \rightarrow \infty$ converges to δ -distribution with PDF concentration at MLE for Bernoulli.



Bayesian approach

Use beta-distribution as a *prior* distribution for our parameter w . From general considerations, the distribution should be symmetrical (unless we have more information):

$$p(w) \sim B(\alpha, \beta).$$

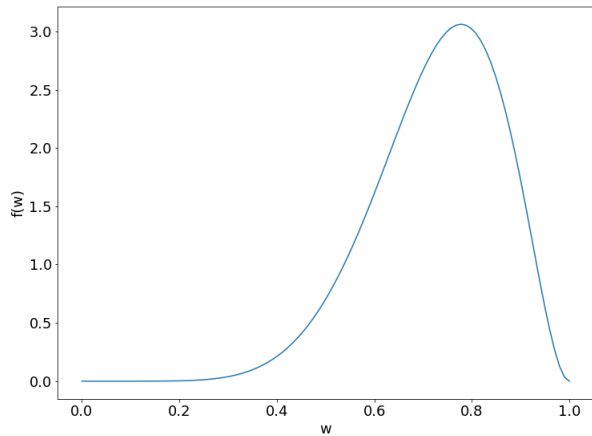
Find the *posterior* distribution of w using Bayes formula:

$$p(w|\mathbf{X}) = \frac{p(\mathbf{X}|w)p(w)}{p(\mathbf{X})} \propto p(\mathbf{X}|w)p(w);$$

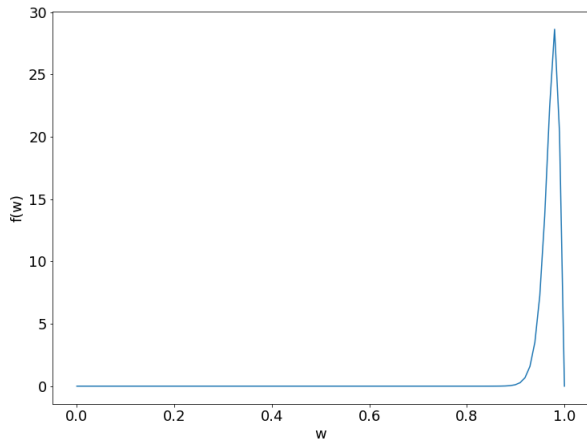
$$\log p(w|\mathbf{X}) = \log p(\mathbf{X}|w) + \log p(w) + \text{Const.}$$

Conclusion: roughly prior is a *regularizer*.

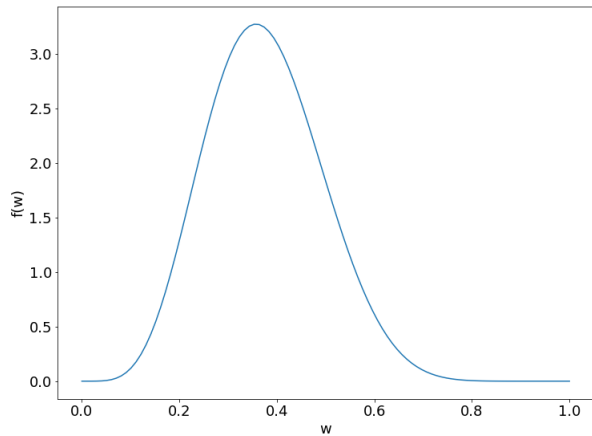
Posterior, $\alpha = 3, \beta = 3$



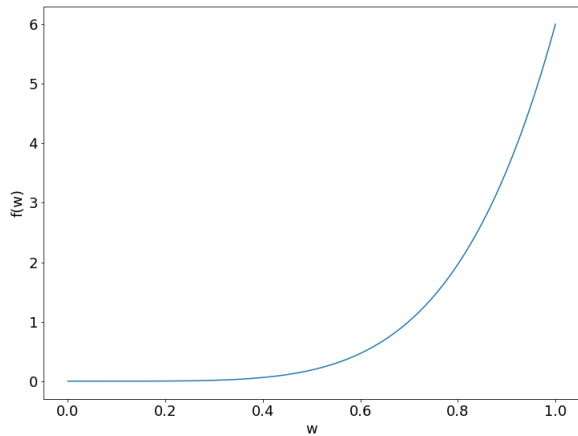
Posterior, 100 elements



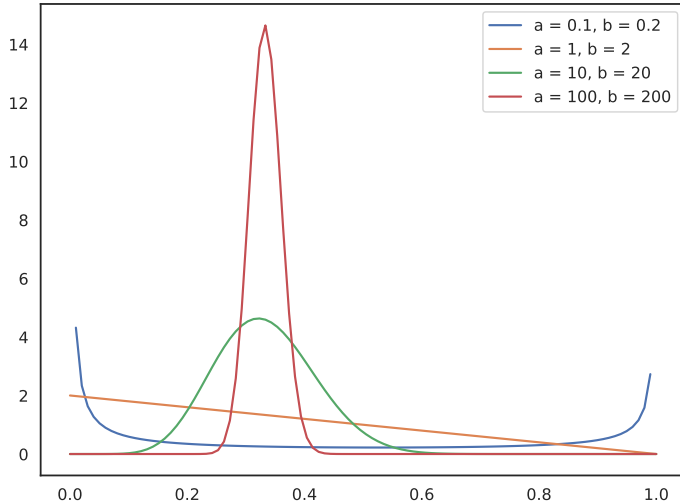
Posterior, $\alpha = 1, \beta = 10$



Posterior, $\alpha = 1, \beta = 1$



Beta-distribution for the sample α - β -ratio



Bayesian inference: first level

Given:

- likelihood $p(\mathbf{X}|\mathbf{w})$ of the dataset \mathbf{X} w.r.t. parameters \mathbf{w} ;
- prior distribution $p(\mathbf{w}|\mathbf{h})$
- prior parameters \mathbf{h} (for the coin problem: $\mathbf{h} = [\alpha, \beta];$)

Then the posterior for \mathbf{w} w.r.t. \mathbf{X} :

$$p(\mathbf{w}|\mathbf{X}, \mathbf{h}) = \frac{p(\mathbf{X}|\mathbf{w})p(\mathbf{w}|\mathbf{h})}{p(\mathbf{X}|\mathbf{h})} \propto p(\mathbf{X}|\mathbf{w})p(\mathbf{w}|\mathbf{h}).$$

Find a point estimate as a maximum posterior probability (MAP):

$$\hat{\mathbf{w}} = \arg \max p(\mathbf{X}|\mathbf{w})p(\mathbf{w}|\mathbf{h}).$$

MAP-estimation is similar to MLE, if

- the dataset is large;
- prior is uniform in an infinitely large region (improper prior)

Why we used Beta-distribution?

$$\begin{aligned} p(w|\mathbf{X}, \alpha, \beta) &\propto p(\mathbf{X}|w)p(w|\alpha, \beta) \propto \\ &\propto w^{\sum x}(1-w)^{m-\sum x} \times w^{\alpha-1}(1-w)^{\beta-1} = \\ &= w^{\alpha-1+\sum x}(1-w)^{m+\beta-\sum x-1} \sim B(\alpha + \sum x, \beta + m - \sum x). \end{aligned}$$

The distribution family is conjugate prior to the likelihood distribution, if the posterior belongs to the same family.

Prior families

- Discrete (labels, discrete parameters)

- ▶ Bernoulli
- ▶ Categorical distributions

Hyperparameters (parameters of the prior parameters):

- ▶ $w \sim \text{Bin}(w)$: $w \sim B(\alpha, \beta)$: conjugate
- ▶ $w \sim \text{Cat}(w)$: $w \sim \text{Dir}(\alpha)$: conjugate

- Real-valued distributions

- ▶ \mathcal{N}
- ▶ Laplace
- ▶ \mathcal{C}

Hyperparameters:

- ▶ Variance, $w \sim \mathcal{N}(\mu, \sigma^2)$, $\sigma^2 \in \Gamma$: conjugate for Gaussian distribution
- ▶ Expectation, $\mu \in \mathcal{N}$: conjugate for Gaussian distribution

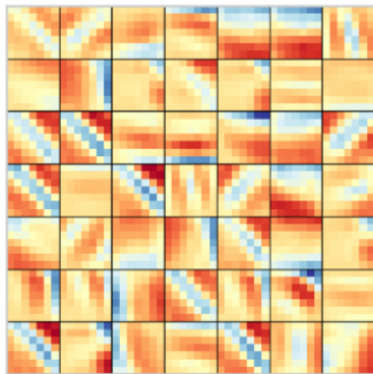
Informative prior vs Uninformative prior

- Informative prior: corresponds to some expert knowledge
 - ▶ Example: air temperature in some region: Gaussian variable with known mean and variance estimated from previous observations.
 - ▶ Mistake in informative prior estimation leads to poor models.
- Uninformative prior: corresponds to some basic knowledge
 - ▶ Example: air temperature in some region: uniform improper prior.
- Weakly-informative prior: somewhere in between
 - ▶ Example: air temperature in some region: uniform distribution in $[-50, 50]$ degrees.

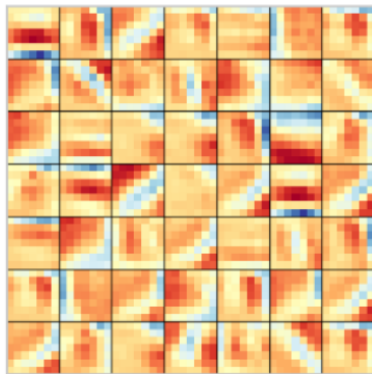
To discuss:

- $\mathbf{w} \sim \mathcal{N}(0, \mathbf{A}^{-1})$ — what type of the prior distribution?
- What if our prior and posterior are very close

The deep weight prior: Atanov et al., 2019



(b) Learned filters



(c) Samples from DWP

The distribution can be modeled implicitly by complex models and can generate rather informative samples.

Jeffreys prior

Uninformative prior:

$$p(\mathbf{w}) \propto \sqrt{\det I(\mathbf{w})} = \sqrt{\det \left(-\frac{\partial^2}{\partial \mathbf{w}^2} \log L(\mathbf{w}) \right)}.$$

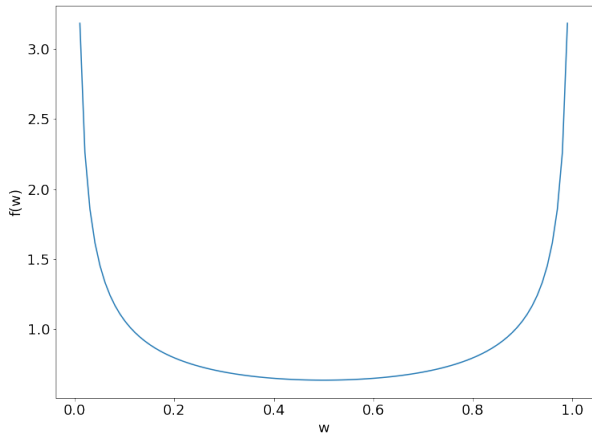
- Invariant under the variable change:

$$p(g(\mathbf{w})) = p(\mathbf{w}) \left| \frac{dg}{d\mathbf{w}} \right| \rightarrow$$

$$p(g(\mathbf{w})) \propto \sqrt{\det I(g(\mathbf{w}))}.$$

- Interpretation: a value inverse to the amount of information obtained by our model from the dataset
- Examples:
 - ▶ $y \in \text{Bin}(w) : p(w) \propto \frac{1}{\sqrt{p(1-p)}} - \text{Beta-distribution } (0.5, 0.5).$
 - ▶ $w \in \mathcal{N}(\mu, \sigma) : p(\mu) \propto \text{Const.}$
 - ▶ $w \in \mathcal{N}(\mu, \sigma) : p(\sigma) \propto \frac{1}{|\sigma|}.$

Uninformative prior \neq flat prior!



See also the talk of Sergey Skorik, 2022

Model selection problem: Bayesian coherent inference

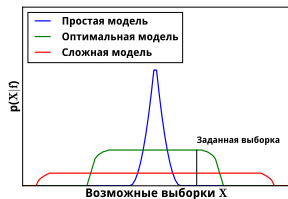
First level: find optimal parameters:

$$\mathbf{w} = \arg \max \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})}{p(\mathcal{D}|\mathbf{h})},$$

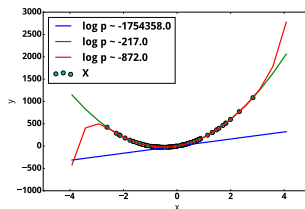
Second level: find model, that give optimal Evidence.

“Evidence”:

$$p(\mathcal{D}|\mathbf{h}) = \int_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})p(\mathbf{w}|\mathbf{h})d\mathbf{w}.$$



Model selection scheme



Example: polynomial regression

Example: linear regression

Given m objects with n features

$\mathbf{f}(\mathbf{X}, \mathbf{w}) = \mathbf{X}\mathbf{w}$; $\mathbf{y} \sim \mathcal{N}(\mathbf{f}(\mathbf{X}, \mathbf{w}), \beta^{-1})$, $\mathbf{w} \sim \mathcal{N}(0, \mathbf{A}^{-1})$.

Write down the integral:

$$\begin{aligned} p(\mathcal{D}|\mathbf{h}) = p(\mathbf{y}|\mathbf{X}, \mathbf{A}, \beta) &= \frac{\sqrt{\beta \cdot |\mathbf{A}|}}{\sqrt{(2\pi)^{m+n}}} \int_{\mathbf{w}} \exp(-0.5\beta(\mathbf{y} - \mathbf{f})^T(\mathbf{y} - \mathbf{f})) \exp(-0.5\mathbf{w}^T \mathbf{A} \mathbf{w}) d\mathbf{w} = \\ &= \frac{\sqrt{\beta \cdot |\mathbf{A}|}}{\sqrt{(2\pi)^{m+n}}} \int_{\mathbf{w}} \exp(-S(\mathbf{w})) d\mathbf{w} \end{aligned}$$

Its value is tractable for the linear regression case:

$$\int_{\mathbf{w}} \exp(-S(\mathbf{w})) d\mathbf{w} = (2\pi)^{\frac{n}{2}} \exp(-S(\hat{\mathbf{w}})) |\mathbf{H}^{-1}|^{0.5},$$

where

$$\mathbf{H} = \mathbf{A} + \beta \mathbf{X}^T \mathbf{X},$$

$$\hat{\mathbf{w}} = \beta \mathbf{H}^{-1} \mathbf{X}^T \mathbf{y}$$

Conclusion: we can find the value of the Evidence for the linear models.

Example: Laplace approximation

Given m objects with n features

$$\mathbf{y} \sim \mathcal{N}(\mathbf{f}(\mathbf{X}, \mathbf{w}), \beta^{-1}), \mathbf{w} \sim \mathcal{N}(0, \mathbf{A}^{-1}).$$

Write down the integral:

$$p(\mathcal{D}|\mathbf{h}) = p(\mathbf{y}|\mathbf{X}, \mathbf{A}, \beta) = \frac{\sqrt{\beta \cdot |\mathbf{A}|}}{\sqrt{(2\pi)^{m+n}}} \int_{\mathbf{w}} \exp(-S(\mathbf{w})) d\mathbf{w}.$$

Use Taylor series for S :

$$S(\mathbf{w}) \approx S(\hat{\mathbf{w}}) + \frac{1}{2} \Delta \mathbf{w}^T \mathbf{H} \Delta \mathbf{w}$$

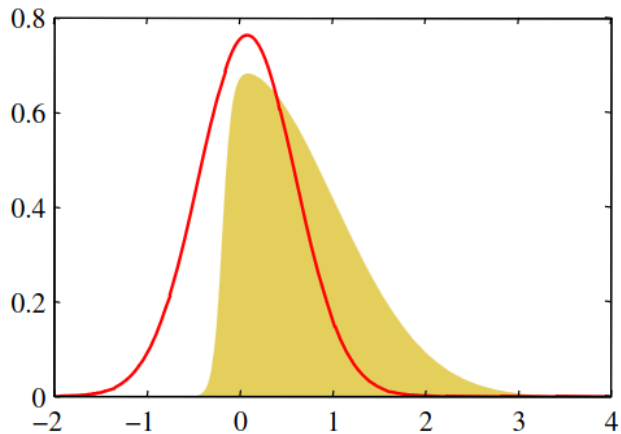
Then:

$$\frac{\sqrt{\beta \cdot |\mathbf{A}|}}{\sqrt{(2\pi)^{m+n}}} S(\hat{\mathbf{w}}) \int_{\mathbf{w}} \exp(-\frac{1}{2} \Delta \mathbf{w}^T \mathbf{H} \Delta \mathbf{w}) d\mathbf{w}$$

The expression corresponds to the PDF for unnormalized Gaussian distribution.

Conclusion: we can use Laplace approximation for the non-linear models.

Laplace approximation: example



Bishop, 2006

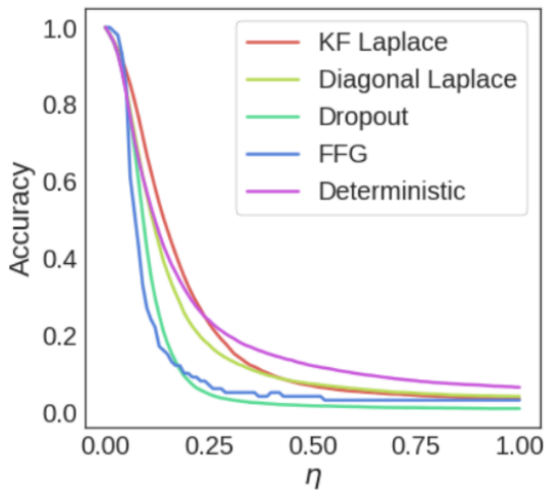
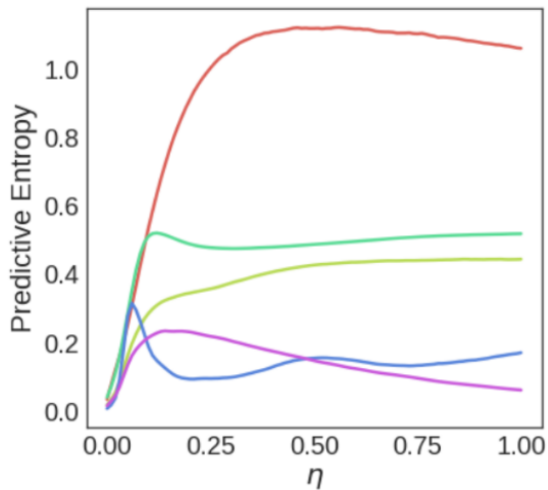
Laplace approximation: drawbacks

- Only Gaussian distribution is available
 - ▶ No multimodality
- Hessian inversion: terribly slow
 - ▶ we can use diagonal matrix, but with worse approximation

A scalable Laplace approximation for neural networks: Ritter et al., 2018

- Decompose the neural network parameters by the layers, make an assumption that parameters from different layers are not correlated
- $\mathbf{H}_l = (\mathbf{f}_l(\mathbf{h}_l)\mathbf{f}_l(\mathbf{h}_l)^\top) \circ \mathbf{H}(\mathbf{h}_l)$ with Kronecker product.
- Reduce the complexity from d^4 to d^2
- Inverse of Kronecker product is equal to the Kronecker product of the inverses

Approximation mode matters



References

- Bishop C. M. Pattern recognition //Machine learning. – 2006. – T. 128. – №. 9.
- MacKay D. J. C., Mac Kay D. J. C. Information theory, inference and learning algorithms. – Cambridge university press, 2003.
- Kuznetsov M., Tokmakova A., Strijov V. Analytic and stochastic methods of structure parameter estimation //Informatica. – 2016. – T. 27. – №. 3. – C. 607-624.
- Coin example: <https://towardsdatascience.com/visualizing-beta-distribution-7391c18031f1>
- Jefreys distribution: <https://medium.datadriveninvestor.com/firths-logistic-regression-classification-with-datasets-that-are-small-imbalanced-or-separated-49d7782a13f1>
- Atanov A. et al. The deep weight prior //arXiv preprint arXiv:1810.06943. – 2018.