# STATS271/371: Applied Bayesian Statistics

**Bayesian Linear Regression**
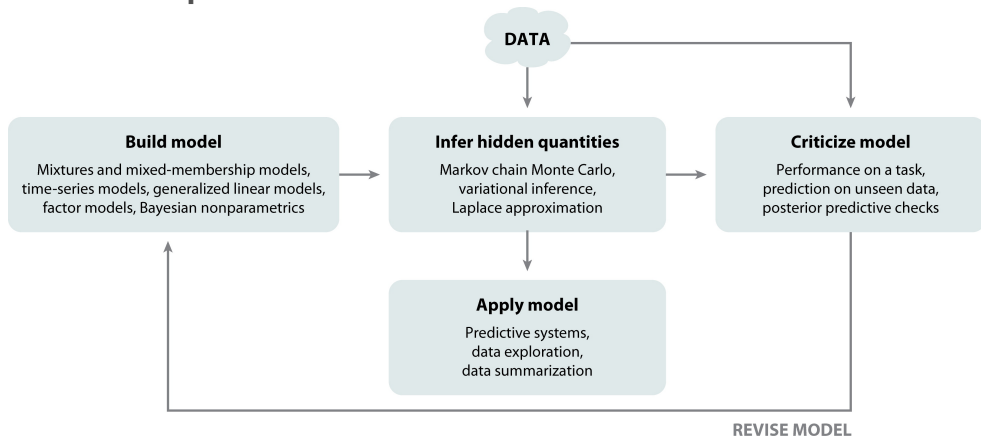
Scott Linderman

March 31, 2021

## Announcements

- ► Please take this short survey: `https://forms.gle/urNME6zwgt1e78Rv6`
- ► Lecture slides are available on Canvas.
- ► Homework 1 will be posted on Friday (Apr 2) and due next Friday (Apr 9).
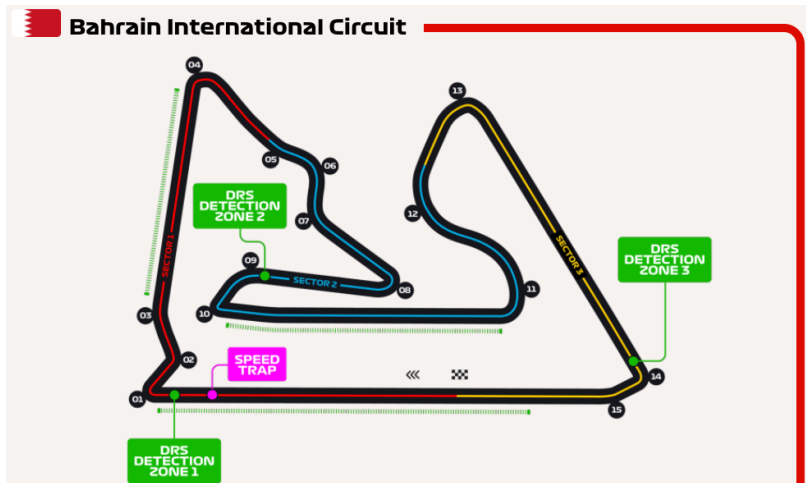
# Lap 1 of Box's Loop



**DATA**

**Build model**

Mixtures and mixed-membership models,
time-series models, generalized linear models,
factor models, Bayesian nonparametrics

**Infer hidden quantities**

Markov chain Monte Carlo,
variational inference,
Laplace approximation

**Criticize model**

Performance on a task,
prediction on unseen data,
posterior predictive checks

**Apply model**

Predictive systems,
data exploration,
data summarization

**REVISE MODEL**

Blei DM. 2014.
Annu. Rev. Stat. Appl. 1:203–32

# Lap 1 of Box's Loop



https://www.formula1.com/en/racing/2021/Bahrain/Circuit.html

## Bayesian Linear Regression

Our first lap around Box's loop will introduce:

▶ **Model:** Bayesian linear regression

▶ **Algorithm:** Exact posterior inference with conjugate priors

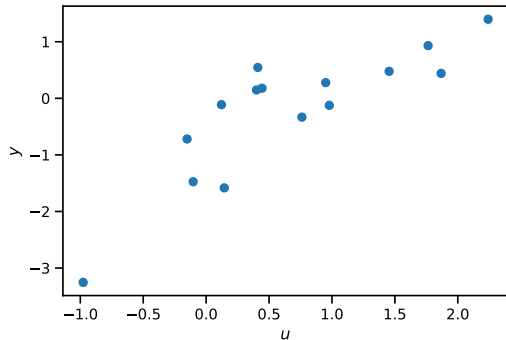▶ **Criticism:** Bayesian model comparison

# Notation

Let

- $y_n \in \mathbb{R}$ denote the $n$-th *observation*
- $\boldsymbol{x}_n \in \mathbb{R}^P$ denote a the *covariates* (aka features) correspond the $n$-th datapoint
- $\boldsymbol{w} \in \mathbb{R}^P$ denote the *weights* of the model
- $\sigma^2 \in \mathbb{R}_+$ denote the variance of the observations

## Example: Polynomial Regression

▶ For example, consider approximating a 1D function $y(u) : \mathbb{R} \to \mathbb{R}$ given noisy observations $\{y_n, u_n\}_{n=1}^N$.

▶ A priori, we don't know if the function is constant, linear, quadratic, cubic, etc.

▶ To fit a polynomial regression model of degree $P-1$, we can encode the inputs $u_n$ with feature vectors,

$$\mathbf{x}_n = (u_n^0, u_n^1, \ldots, u_n^{P-1}) \in \mathbb{R}^P \qquad (1)$$

and perform a linear regression.

## Likelihood

We assume a standard Gaussian likelihood with independent noise for each datapoint,

$$p(\{y_n\}_{n=1}^N \mid \{\boldsymbol{x}_n\}_{n=1}^N, \boldsymbol{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(y_n \mid \boldsymbol{w}^\top \boldsymbol{x}_n, \sigma^2) \tag{2}$$

$$= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_n - \boldsymbol{w}^\top \boldsymbol{x}_n)^2\right\} \tag{3}$$

$$= \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2}\frac{y_n^2}{\sigma^2} + \frac{y_n \boldsymbol{x}_n^\top \boldsymbol{w}}{\sigma^2} - \frac{1}{2}\frac{\boldsymbol{w}^\top \boldsymbol{x}_n \boldsymbol{x}_n^\top \boldsymbol{w}}{\sigma^2}\right\} \tag{4}$$

$$\propto (\sigma^2)^{-\frac{N}{2}} \exp\left\{-\frac{1}{2}\Big\langle \sum_{n=1}^N y_n^2, \frac{1}{\sigma^2} \Big\rangle + \Big\langle \sum_{n=1}^N y_n \boldsymbol{x}_n, \frac{\boldsymbol{w}}{\sigma^2} \Big\rangle - \frac{1}{2}\Big\langle \sum_{n=1}^N \boldsymbol{x}_n \boldsymbol{x}_n^\top, \frac{\boldsymbol{w}\boldsymbol{w}^\top}{\sigma^2} \Big\rangle\right\} \tag{5}$$

The *sufficient statistics* of the data are $\left(\sum_{n=1}^N y_n^2, \sum_{n=1}^N y_n \boldsymbol{x}_n, \sum_{n=1}^N \boldsymbol{x}_n \boldsymbol{x}_n^\top\right)$.

## Aside: inner products between two matrices

The following equalities hold for the scalar quadratic form above,

$$w^\top x_n x_n^\top w = \mathrm{Tr}(w^\top x_n x_n^\top w) \tag{6}$$

$$= \mathrm{Tr}(x_n x_n^\top w w^\top) \tag{7}$$

$$= \sum_{i=1}^{P} \sum_{j=1}^{P} [x_n x_n^\top]_{ij} [w w^\top]_{ji} \tag{8}$$

$$\triangleq \langle x_n x_n^\top, w w^\top \rangle. \tag{9}$$

The inner product between two matrices $x_n x_n^\top$ and $w w^\top$ is defined by the last expression. As the sum of the element-wise product, it naturally generalizes the inner product between two vectors.

## Review of maximum likelihood estimation

Before considering a Bayesian treatment, let's recall the standard maximum likelihood estimate of the parameters.

The log likelihood is,

$$\mathcal{L}(\mathbf{w}, \sigma^2) = \log p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N, \mathbf{w}, \sigma^2) \tag{10}$$

$$= -\frac{N}{2} \log \sigma^2 - \frac{1}{2} \Big\langle \sum_{n=1}^N y_n^2, \frac{1}{\sigma^2} \Big\rangle + \Big\langle \sum_{n=1}^N y_n \mathbf{x}_n, \frac{\mathbf{w}}{\sigma^2} \Big\rangle - \frac{1}{2} \Big\langle \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top, \frac{\mathbf{w}\mathbf{w}^\top}{\sigma^2} \Big\rangle \tag{11}$$

## Review of maximum likelihood estimation

Taking the gradient and setting it to zero,

$$\nabla_{\boldsymbol{w}}\mathscr{L}(\boldsymbol{w}, \sigma^2) = \frac{1}{\sigma^2}\sum_{n=1}^{N}y_n\boldsymbol{x}_n - \left(\frac{1}{\sigma^2}\sum_{n=1}^{N}\boldsymbol{x}_n\boldsymbol{x}_n^\top\right)\boldsymbol{w} = 0 \tag{12}$$

$$\implies \boldsymbol{w}_{\text{MLE}} = \left(\sum_{n=1}^{N}\boldsymbol{x}_n\boldsymbol{x}_n^\top\right)^{-1}\left(\sum_{n=1}^{N}y_n\boldsymbol{x}_n\right). \tag{13}$$

Letting

$$\boldsymbol{X} = \begin{bmatrix} - & \boldsymbol{x}_1^\top & - \\ & \vdots & \\ - & \boldsymbol{x}_N^\top & - \end{bmatrix}, \quad \text{and} \quad \boldsymbol{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \tag{14}$$

we can write this in the more familiar form of the ordinary least squares solution,

$$\boldsymbol{w}_{\text{MLE}} = (\boldsymbol{X}^\top\boldsymbol{X})^{-1}\boldsymbol{X}^\top\boldsymbol{y}. \tag{15}$$

## Review of maximum likelihood estimation II

Now let $\hat{y} = Xw_{\text{MLE}} = X(X^\top X)^{-1}X^\top y$ denote the predicted observations under the optimal weights. Substituting this in, we have

$$\mathscr{L}(w_{\text{MLE}}, \sigma^2) = -\frac{N}{2}\log\sigma^2 - \frac{1}{2\sigma^2}(y - \hat{y})^\top(y - \hat{y}) \tag{16}$$

Taking derivatives wrt $1/\sigma^2$ and setting to zero,

$$\frac{\partial}{\partial\sigma^{-2}}\mathscr{L}(w_{\text{MLE}}, \sigma^2) = \frac{N}{2}\sigma^2 - \frac{1}{2}(y - \hat{y})^\top(y - \hat{y}) = 0 \tag{17}$$

$$\implies \sigma^2_{\text{MLE}} = \frac{1}{N}(y - \hat{y})^\top(y - \hat{y}) \tag{18}$$

$$= \frac{1}{N}(y^\top y - y^\top X(X^\top X)^{-1}X^\top y) \tag{19}$$

$$= \frac{1}{N}y^\top(I - H)y, \tag{20}$$

where $H = X(X^\top X)^{-1}X^\top$ is the *hat matrix*, which projects onto the span of the columns of $X$.

## Prior

Now consider a Bayesian treatment in which we introduce a prior on the parameters $w$ and $\sigma^2$. Some desiderata when choosing a prior:

▶ It should capture intuition about the weights, like the general scale or sparsity.

▶ In the case where we have little prior information, it should be a broad and relatively uninformative distribution.

▶ All else equal, we'd prefer if it permits tractable posterior calculations.

## Prior II

Let's assume we don't know much about the weights *a priori*. We'll choose a prior of the following form,

$$p(\boldsymbol{w}, \sigma^2) = \mathrm{Inv}\text{-}\chi^2(\sigma^2 \mid \nu, \tau^2) \, \mathcal{N}(\boldsymbol{w} \mid \boldsymbol{\mu}, \sigma^2 \boldsymbol{\Lambda}^{-1}), \tag{21}$$

where

▶ $\nu, \tau^2 \in \mathbb{R}_+$ are the degrees-of-freedom and scaling parameter, respectively, of the *inverse chi-squared distribution*.

▶ $\boldsymbol{\mu} \in \mathbb{R}^P$ and $\boldsymbol{\Lambda} \in \mathbb{R}_{>0}^{P \times P}$ are the mean and (positive definite) precision matrix, respectively, of a *multivariate normal distribution*.

## Aside: Inverse Chi-Squared Distribution

[From Wikipedia] Let $s^2$ be the sample mean of the squares of $\nu$ independent normal random variables with mean 0 and precision $\tau^2$. Then $\sigma^2 = 1/s^2$ is distributed as $\text{Inv}{-}\chi^2(\nu, \tau^2)$ and has pdf,
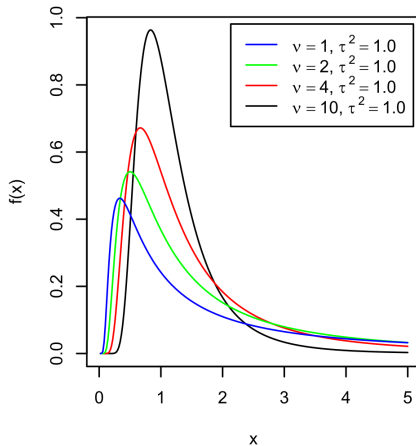
$$\text{Inv}{-}\chi^2(\sigma^2 \mid \nu, \tau^2) = \frac{\left(\frac{\tau^2 \nu}{2}\right)^{\frac{\nu}{2}}}{\Gamma(\frac{\nu}{2})} (\sigma^2)^{-(1+\frac{\nu}{2})} \exp\left\{-\frac{1}{2}\left\langle \nu\tau^2, \frac{1}{\sigma^2} \right\rangle\right\} \tag{22}$$

The scaled inverse chi-squared distribution is a reparametrization of the inverse gamma distribution. Specifically, if

$$\sigma^2 \sim \text{Inv}{-}\chi^2(\nu, \tau^2) \quad \Longleftrightarrow \quad \sigma^2 \sim \text{IGa}\left(\frac{\nu}{2}, \frac{\nu\tau^2}{2}\right). \tag{23}$$

This reparameterization is sometimes easier to work with as a conjugate prior for the variance of a Gaussian distribution.

# Aside Inverse Chi-Squared Distribution II



https://en.wikipedia.org/wiki/Scaled_inverse_chi-squared_distribution

## Prior density

Now expanding the prior density,

$$p(\boldsymbol{w}, \sigma^2) = \frac{\left(\frac{\tau^2 \nu}{2}\right)^{\frac{\nu}{2}}}{\Gamma(\frac{\nu}{2})} (\sigma^2)^{-(1+\frac{\nu}{2})} e^{-\frac{\nu \tau^2}{2\sigma^2}} \times (2\pi)^{-\frac{p}{2}} |\sigma^2 \boldsymbol{\Lambda}^{-1}|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2}(\boldsymbol{w} - \boldsymbol{\mu})^\top \boldsymbol{\Lambda}(\boldsymbol{w} - \boldsymbol{\mu})\right\} \qquad (24)$$

$$= \frac{1}{Z(\nu, \tau^2, \boldsymbol{\Lambda})} (\sigma^2)^{-(1+\frac{\nu}{2}+\frac{p}{2})} \exp\left\{-\frac{1}{2}\left\langle \nu\tau^2 + \boldsymbol{\mu}^\top \boldsymbol{\Lambda}\boldsymbol{\mu}, \frac{1}{\sigma^2}\right\rangle + \left\langle \boldsymbol{\Lambda}\boldsymbol{\mu}, \frac{\boldsymbol{w}}{\sigma^2}\right\rangle - \frac{1}{2}\left\langle \boldsymbol{\Lambda}, \frac{\boldsymbol{w}\boldsymbol{w}^\top}{\sigma^2}\right\rangle\right\} \qquad (25)$$
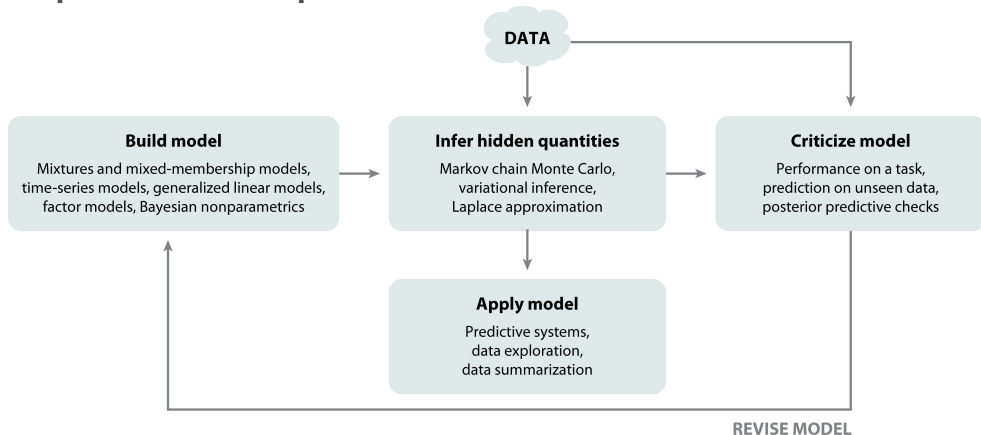
where the normalizing constant is

$$Z(\nu, \tau^2, \boldsymbol{\Lambda}) = \frac{\Gamma(\frac{\nu}{2})}{\left(\frac{\tau^2 \nu}{2}\right)^{\frac{\nu}{2}}} (2\pi)^{\frac{p}{2}} |\boldsymbol{\Lambda}|^{-\frac{1}{2}} \qquad (26)$$

## Properties of the prior

► **Conjugacy:** Note that the functional form is the same as in the likelihood, (5).

  ► In the exponent, both are linear functions of $\frac{1}{\sigma^2}$, $\frac{w}{\sigma^2}$, and $\frac{ww^\top}{\sigma^2}$.

  ► This is the defining property of a *conjugate prior*, and it will lead to a closed form posterior distribution.

► **Uninformativeness:** As $\nu, \Lambda \to 0$, the prior reduces to an *improper* prior of the form $p(w, \sigma^2) \propto (\sigma^2)^{-(1+\frac{p}{2})}$.

  ► That is, it is effectively uniform in the weights and shrinking polynomially in the variance.

# Box's Loop: Infer hidden quantities



**DATA**

**Build model**

Mixtures and mixed-membership models, time-series models, generalized linear models, factor models, Bayesian nonparametrics

**Infer hidden quantities**

Markov chain Monte Carlo, variational inference, Laplace approximation

**Criticize model**

Performance on a task, prediction on unseen data, posterior predictive checks

**Apply model**

Predictive systems, data exploration, data summarization

**REVISE MODEL**

Blei DM. 2014.
Annu. Rev. Stat. Appl. 1:203–32

# Algorithm

► Thanks to the conjugacy of the prior, we can perform *exact* posterior inference in this model.

► Note that this is a very special case!

► The remainder of the models we'll encounter in this course will not be so nice, and we'll have to make some approximations to the posterior distribution.

## Posterior Distribution

For this well behaved model we have,

$$p(\boldsymbol{w}, \sigma^2 \mid \{\boldsymbol{x}_n, y_n\}_{n=1}^N) \propto p(\boldsymbol{w}, \sigma^2)\, p(\{y_n\}_{n=1}^N \mid \{\boldsymbol{x}_n\}_{n=1}^N, \boldsymbol{w}, \sigma^2) \tag{27}$$

$$\propto (\sigma^2)^{-(1+\frac{\nu}{2}+\frac{P}{2}+\frac{N}{2})} \times \exp\left\{ -\frac{1}{2}\Big\langle \nu\tau^2 + \boldsymbol{\mu}^\top \boldsymbol{\Lambda}\boldsymbol{\mu} + \sum_{n=1}^N y_n^2, \frac{1}{\sigma^2} \Big\rangle \right.$$

$$\left. + \Big\langle \boldsymbol{\Lambda}\boldsymbol{\mu} + \sum_{n=1}^N y_n \boldsymbol{x}_n, \frac{\boldsymbol{w}}{\sigma^2} \Big\rangle - \frac{1}{2}\Big\langle \boldsymbol{\Lambda} + \sum_{n=1}^N \boldsymbol{x}_n \boldsymbol{x}_n^\top, \frac{\boldsymbol{w}\boldsymbol{w}^\top}{\sigma^2} \Big\rangle \right\}$$

## Posterior Distribution II

Again, we see this is the same family as the prior,

$$p(\boldsymbol{w}, \sigma^2 \mid \{\boldsymbol{x}_n, y_n\}_{n=1}^N) = \text{Inv}-\chi^2(\sigma^2 \mid \nu', \tau'^2) \, \mathcal{N}(\boldsymbol{w} \mid \boldsymbol{\mu}', \sigma^2 \Lambda'^{-1}), \tag{28}$$

where

$$\Lambda' = \Lambda + \sum_{n=1}^N \boldsymbol{x}_n \boldsymbol{x}_n^\top \tag{29}$$

$$\nu' = \nu + N \tag{30}$$

$$\boldsymbol{\mu}' = \Lambda'^{-1} \left( \Lambda \boldsymbol{\mu} + \sum_{n=1}^N y_n \boldsymbol{x}_n \right) \tag{31}$$

$$\tau'^2 = \frac{1}{\nu'} \left( \nu \tau^2 + \boldsymbol{\mu}^\top \Lambda \boldsymbol{\mu} + \sum_{n=1}^N y_n^2 - \boldsymbol{\mu}'^\top \Lambda' \boldsymbol{\mu}' \right) \tag{32}$$

## Uninformative limit

Consider the uninformative limit in which $\nu, \Lambda \to 0$. Then,

$$\Lambda' \to \sum_{n=1}^{N} \boldsymbol{x}_n \boldsymbol{x}_n^\top \tag{33}$$

$$\nu' \to N \tag{34}$$

$$\boldsymbol{\mu}' \to \left( \sum_{n=1}^{N} \boldsymbol{x}_n \boldsymbol{x}_n^\top \right)^{-1} \left( \sum_{n=1}^{N} y_n \boldsymbol{x}_n \right) \tag{35}$$

$$\tau'^2 \to \frac{1}{N} \left( \sum_{n=1}^{N} y_n^2 - \left( \sum_{n=1}^{N} y_n \boldsymbol{x}_n \right)^\top \left( \sum_{n=1}^{N} \boldsymbol{x}_n \boldsymbol{x}_n^\top \right)^{-1} \left( \sum_{n=1}^{N} y_n \boldsymbol{x}_n \right) \right) \tag{36}$$

Or in matrix notation

$$\Lambda' \to \boldsymbol{X}^\top \boldsymbol{X} \qquad\qquad \nu' \to N \tag{37}$$

$$\boldsymbol{\mu}' \to \left( \boldsymbol{X}^\top \boldsymbol{X} \right)^{-1} \left( \boldsymbol{X}^\top \boldsymbol{y} \right) \qquad\qquad \tau'^2 \to \frac{1}{N} \boldsymbol{y}^\top (\boldsymbol{I} - \boldsymbol{H}) \boldsymbol{y} \tag{38}$$

## Posterior Mode (aka MAP Estimate)

Under this uninformative prior, the posterior mode, aka the *maximum a posteriori* (MAP) estimate, is,

$$\boldsymbol{w}_{\text{MAP}} = \left(\boldsymbol{X}^\top \boldsymbol{X}\right)^{-1} \left(\boldsymbol{X}^\top \boldsymbol{y}\right) \tag{39}$$

$$\sigma_{\text{MAP}}^2 = \frac{\nu' \tau'^2}{\nu' + 2} = \frac{\boldsymbol{y}^\top (\boldsymbol{I} - \boldsymbol{H}) \boldsymbol{y}}{N + 2}. \tag{40}$$
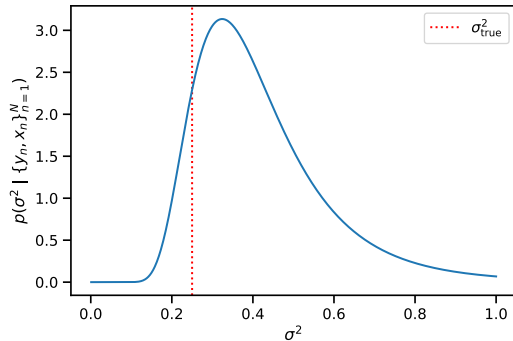
In other words, $\boldsymbol{w}_{\text{MAP}} = \boldsymbol{w}_{\text{MLE}}$ and $\sigma_{\text{MAP}}^2 = \frac{N}{N+2} \sigma_{\text{MLE}}^2$.

The weights are unchanged and the variance is slightly smaller.
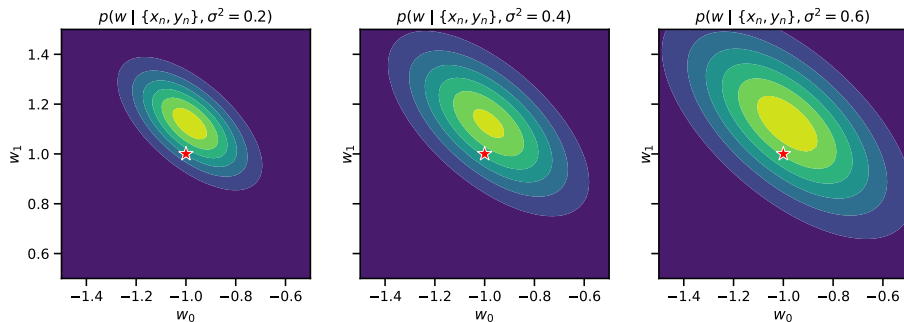
# Posterior distribution for synthetic example

First we plot the posterior distribution of the variance,

$$p(\sigma^2 \mid \{y_n, \boldsymbol{x}_n\}_{n=1}^{N}) = \mathrm{IGa}(\nu', \tau'^2) \qquad (41)$$

# Posterior distribution for synthetic example

Then plot $p(\mathbf{w} \mid \{y_n, \mathbf{x}_n\}_{n=1}^{N}, \sigma^2) = \mathcal{N}(\mathbf{w} \mid \boldsymbol{\mu}', \sigma^2 \boldsymbol{\Lambda}'^{-1})$ for a few values of $\sigma^2$

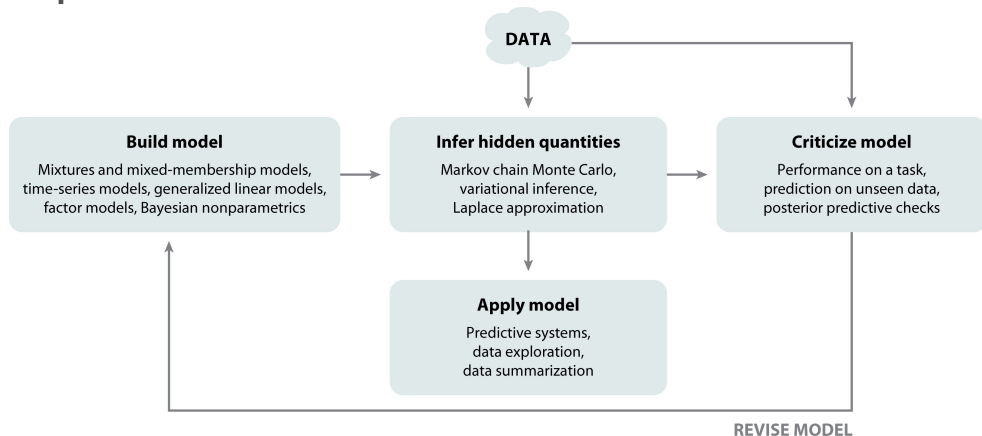## Exercise: $L_2$ regularization

What happens to $w_{\text{MAP}}$ if you set $\Lambda = \lambda I$ for scalar $\lambda > 0$ and set $\mu = 0$?

This is known as $L_2$ regularization, Tikhonov regularization, or "weight decay" in various communities.

# Box's Loop: Criticize model



**DATA**

**Build model**
Mixtures and mixed-membership models, time-series models, generalized linear models, factor models, Bayesian nonparametrics

**Infer hidden quantities**
Markov chain Monte Carlo, variational inference, Laplace approximation

**Criticize model**
Performance on a task, prediction on unseen data, posterior predictive checks

**Apply model**
Predictive systems, data exploration, data summarization

**REVISE MODEL**

Blei DM. 2014.
Annu. Rev. Stat. Appl. 1:203–32

## Model Comparison

▶ The marginal likelihood, aka model evidence, is a useful measure of how well a model fits the data.

▶ Specifically, it measures the *expected* probability assigned to the data, integrating over possible parameters under the prior,

$$p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N) = \int p(\mathbf{w}, \sigma^2) p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N, \mathbf{w}, \sigma^2) \, d\mathbf{w} \, d\sigma^2 \tag{42}$$

$$= \mathbb{E}_{p(\mathbf{w}, \sigma^2)} \left[ p(\{y_n\}_{n=1}^N \mid \{\mathbf{x}_n\}_{n=1}^N, \mathbf{w}, \sigma^2) \right] \tag{43}$$

▶ If a prior distribution puts high probability on weights and variances that then assign high conditional probability to the given data, the marginal likelihood will be large.

▶ If the prior spreads its probability mass over a wide range of weights, it may have a lower marginal likelihood than one that concentrates mass around the weights that achieve maximal likelihood.

## Marginal Likelihood

Under the conjugate prior above, we can compute the marginal likelihood in closed form,

$$
p(\{y_n\}_{n=1}^N \mid \{\boldsymbol{x}_n\}_{n=1}^N) = \int \frac{(2\pi)^{-\frac{N}{2}}}{Z(\nu, \tau^2, \Lambda)} (\sigma^2)^{-(1+\frac{\nu'}{2}+\frac{p}{2})}
$$

$$
\exp\left\{ -\frac{1}{2}\left\langle \nu'\tau'^2 + \boldsymbol{\mu}'^\top \Lambda' \boldsymbol{\mu}', \frac{1}{\sigma^2}\right\rangle \right.
$$

$$
\left. + \left\langle \Lambda'\boldsymbol{\mu}', \frac{\boldsymbol{w}}{\sigma^2}\right\rangle - \frac{1}{2}\left\langle \Lambda', \frac{\boldsymbol{w}\boldsymbol{w}^\top}{\sigma^2}\right\rangle \right\} d\boldsymbol{w}\, d\sigma^2 \tag{44}
$$

$$
= (2\pi)^{-\frac{N}{2}} \frac{Z(\nu', \tau'^2, \Lambda')}{Z(\nu, \tau^2, \Lambda)} \int \frac{1}{Z(\nu', \tau'^2, \Lambda')} \text{``} \qquad \cdots \qquad \text{''} \, d\boldsymbol{w}\, d\sigma^2 \tag{45}
$$

$$
= (2\pi)^{-\frac{N}{2}} \frac{Z(\nu', \tau'^2, \Lambda')}{Z(\nu, \tau, \Lambda)} \tag{46}
$$

## Marginal Likelihood II

Under the conjugate prior above, we can compute the marginal likelihood in closed form,

$$p(\{y_n\}_{n=1}^N \mid \{\boldsymbol{x}_n\}_{n=1}^N) = (2\pi)^{-\frac{N}{2}} \frac{Z(\nu', \tau'^2, \boldsymbol{\Lambda}')}{Z(\nu, \tau, \boldsymbol{\Lambda})} \tag{47}$$
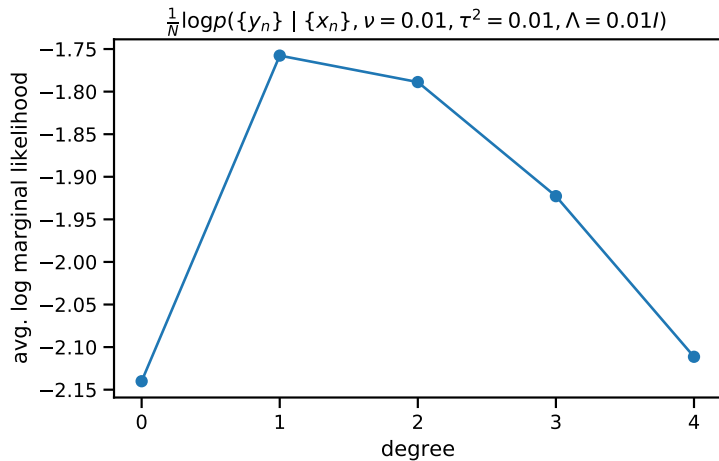
$$= (2\pi)^{-\frac{N}{2}} \frac{\Gamma(\frac{\nu'}{2})}{\Gamma(\frac{\nu}{2})} \frac{(\frac{\tau^2 \nu}{2})^{\frac{\nu}{2}}}{(\frac{\tau'^2 \nu'}{2})^{\frac{\nu'}{2}}} \frac{|\boldsymbol{\Lambda}|^{\frac{1}{2}}}{|\boldsymbol{\Lambda}'|^{\frac{1}{2}}} \tag{48}$$

## Properly speaking...

Note that in order for the marginal likelihood to be meaningful, we need to have a *proper* prior distribution.

In the uninformative/improper limit, the marginal likelihood goes to zero.

# Example: Using the marginal likelihood to select degree of a polynomial regression



Figure: plot with title $\frac{1}{N}\log p(\{y_n\} \mid \{x_n\}, \nu = 0.01, \tau^2 = 0.01, \Lambda = 0.01 I)$, x-axis "degree", y-axis "avg. log marginal likelihood".

## Exercise: Unpacking the marginal likelihood

Consider selecting the degree of a polynomial regression by maximizing the marginal likelihood above. Which ratios in the marginal likelihood are growing, shrinking, or fixed, as you increase the degree *P*?

# Preview: Posterior Predictive Distribution

▶ One of the main uses of regression models is to make predictions, e.g. of $y_{N+1}$ at $\boldsymbol{x}_{N+1}$.

▶ In Bayesian data analysis, this is given by the *posterior predictive distribution*,

$$p(y_{N+1} \mid \boldsymbol{x}_{N+1}, \{y_n, \boldsymbol{x}_n\})_{n=1}^N) = \int p(y_{N+1} \mid \boldsymbol{x}_{N+1}, \boldsymbol{w}, \sigma^2)\, p(\boldsymbol{w}, \sigma^2 \mid \{y_n, \boldsymbol{x}_n\}_{n=1}^N\, \mathrm{d}\boldsymbol{w}\, \mathrm{d}\sigma^2 \quad (49)$$

▶ Generally, we can approximate the posterior predictive distribution with Monte Carlo.

▶ For Bayesian linear regression with a conjugate prior, we can compute it in closed form.

## Preview: Posterior Predictive Distribution II

We have,

$$p(y_{N+1} \mid \mathbf{x}_{N+1}, \{y_n, \mathbf{x}_n\})_{n=1}^N) = \int p(y_{N+1} \mid \mathbf{x}_{N+1}, \mathbf{w}, \sigma^2) \, p(\mathbf{w}, \sigma^2 \mid \{y_n, \mathbf{x}_n\}_{n=1}^N \, \mathrm{d}\mathbf{w} \, \mathrm{d}\sigma^2 \tag{50}$$
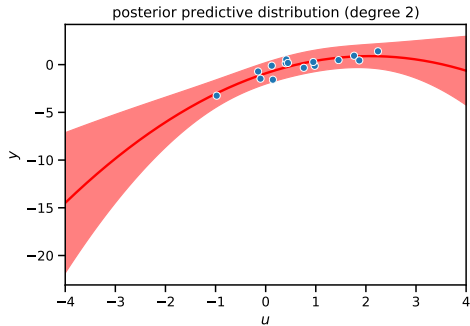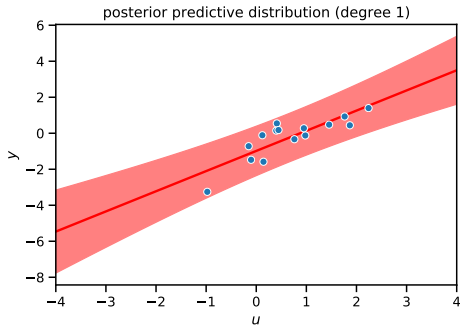
$$= \int \mathcal{N}(y_{N+1} \mid \mathbf{w}^\top \mathbf{x}_{N+1}, \sigma^2) \, \mathcal{N}(\mathbf{w} \mid \boldsymbol{\mu}', \sigma^2 \mathbf{\Lambda}'^{-1}) \, \mathrm{Inv}{-}\chi^2(\sigma^2 \mid \nu', \tau'^2) \, \mathrm{d}\mathbf{w} \, \mathrm{d}\sigma^2 \tag{51}$$

$$= \int \mathcal{N}(y_{N+1} \mid \boldsymbol{\mu}'^\top \mathbf{x}_{N+1}, \sigma^2(1 + \mathbf{x}_{N+1}^\top \mathbf{\Lambda}^{-1} \mathbf{x}_{N+1})) \, \mathrm{Inv}{-}\chi^2(\sigma^2 \mid \nu', \tau'^2) \, \mathrm{d}\sigma^2 \tag{52}$$

$$= t(y_{N+1} \mid \nu', \boldsymbol{\mu}'^\top \mathbf{x}_{N+1}, \tau'^2(1 + \mathbf{x}_{N+1}^\top \mathbf{\Lambda}^{-1} \mathbf{x}_{N+1})) \tag{53}$$

where $t(\cdot \mid \nu, \mu, \tau^2)$ is the density of a (generalized) *Students-t* distribution with $\nu$ degrees of freedom, location $\mu$, and scale $\tau$.

# Preview: Posterior predictive distribution III

## Bonus: Multivariate observations

Now consider multivariate observations $\boldsymbol{y}_n \in \mathbb{R}^D$ and a likelihood,

$$p(\{\boldsymbol{y}_n\}_{n=1}^N \mid \{\boldsymbol{x}_n\}_{n=1}^N, \boldsymbol{W}, \boldsymbol{S}) = \prod_{n=1}^N \mathcal{N}(\boldsymbol{y}_n \mid \boldsymbol{W}\boldsymbol{x}_n, \boldsymbol{S}) \tag{54}$$

where $\boldsymbol{W} \in \mathbb{R}^{D \times P}$ is now a weight *matrix* and $\boldsymbol{S} \in \mathbb{R}_{\succ 0}^{D \times D}$ is a positive definite covariance matrix.

## Bonus: Multivariate observations II

Expanding the likelihood, as above, we obtain,

$$p(\{\boldsymbol{y}_n\}_{n=1}^N \mid \{\boldsymbol{x}_n\}_{n=1}^N, \boldsymbol{W}, \boldsymbol{S}) \propto |\boldsymbol{S}|^{-\frac{N}{2}} \exp\left\{ -\frac{1}{2}\left\langle \sum_{n=1}^N \boldsymbol{y}_n\boldsymbol{y}_n^\top, \boldsymbol{S}^{-1} \right\rangle \right.$$
$$+ \left\langle \sum_{n=1}^N \boldsymbol{y}_n\boldsymbol{x}_n^\top, \boldsymbol{S}^{-1}\boldsymbol{W} \right\rangle$$
$$\left. -\frac{1}{2}\left\langle \sum_{n=1}^N \boldsymbol{x}_n\boldsymbol{x}_n^\top, \boldsymbol{W}^\top\boldsymbol{S}^{-1}\boldsymbol{W} \right\rangle \right\} \tag{55}$$

# Conjugate prior

This is conjugate with a *matrix normal inverse Wishart* (MNIW) prior of the form,

$$p(\boldsymbol{W}, \boldsymbol{S}) = \mathrm{MNIW}(\boldsymbol{W}, \boldsymbol{S} \mid \boldsymbol{\Psi}, \nu, \boldsymbol{\mu}, \boldsymbol{\Lambda}) \tag{56}$$

$$= \mathrm{IW}(\boldsymbol{S} \mid \boldsymbol{\Psi}, \nu) \, \mathrm{MN}(\boldsymbol{W} \mid \boldsymbol{\mu}, \boldsymbol{S}, \boldsymbol{\Lambda}^{-1}). \tag{57}$$

## Inverse Wishart Distribution

The first the first term is an *inverse Wishart* distribution,

$$\mathrm{IW}(\boldsymbol{S} \mid \boldsymbol{\Psi}, \nu) = \frac{\left(\frac{|\boldsymbol{\Psi}|}{2^D}\right)^{\frac{\nu}{2}}}{\Gamma_D\left(\frac{\nu}{2}\right)} |\boldsymbol{S}|^{-(\nu+D+1)/2} \exp\left\{-\frac{1}{2}\langle\boldsymbol{\Psi}, \boldsymbol{S}^{-1}\rangle\right\} \tag{58}$$

where $\Gamma_D(\cdot)$ denotes the multivariate gamma function. The inverse Wishart is a multivariate generalization of the scaled inverse chi-squared distribution.

## Matrix Normal Distribution

The second term is a *matrix normal* distribution,

$$\mathrm{MN}(\boldsymbol{W} \mid \boldsymbol{\mu}, \boldsymbol{S}, \Lambda^{-1}) = \mathcal{N}(\mathrm{vec}(\boldsymbol{W}) \mid \mathrm{vec}(\boldsymbol{\mu}), \boldsymbol{S} \otimes \Lambda^{-1}) \tag{59}$$

$$= (2\pi)^{-\frac{DP}{2}} |\boldsymbol{S}|^{-\frac{P}{2}} |\Lambda|^{\frac{D}{2}} \exp\left\{-\frac{1}{2}\mathrm{Tr}\left(\Lambda(\boldsymbol{W}-\boldsymbol{\mu})^\top \boldsymbol{S}^{-1}(\boldsymbol{W}-\boldsymbol{\mu})\right)\right\} \tag{60}$$

$$= (2\pi)^{-\frac{DP}{2}} |\boldsymbol{S}|^{-\frac{P}{2}} |\Lambda|^{\frac{D}{2}} \exp\left\{-\frac{1}{2}\left\langle \boldsymbol{\mu}\Lambda\boldsymbol{\mu}^\top, \boldsymbol{S}^{-1}\right\rangle \right.$$
$$\left. + \left\langle \boldsymbol{\mu}\Lambda, \boldsymbol{S}^{-1}\boldsymbol{W}\right\rangle \right.$$
$$\left. -\frac{1}{2}\left\langle \Lambda, \boldsymbol{W}^\top \boldsymbol{S}^{-1}\boldsymbol{W}\right\rangle \right\} \tag{61}$$

The product of the matrix normal and inverse Wishart densities has natural parameters $\log|\boldsymbol{S}|$, $\boldsymbol{S}^{-1}$, $\boldsymbol{S}^{-1}\boldsymbol{W}$, and $\boldsymbol{W}^\top \boldsymbol{S}^{-1}\boldsymbol{W}$.

# Exercise: Matrix Normal Inverse Wishart Distribution

*Show that the prior used for the scalar observations is a special case of the MNIW prior.*