

Machine Learning

Supervised Machine Learning – Regression

Part 4: ML and Bayesian regression

The “Bias” Term

Data Model Reminder

- Up until now:

$$y = f(\mathbf{x}) + \epsilon$$

For parametric regression, we have some idea on the structure of f (encoded with parameters)

$$\text{E.g., } f(\mathbf{x}) = \mathbf{b}(\mathbf{x})^T \mathbf{w}, \text{ for } \mathbf{b}(\mathbf{x}) = [1, \phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})]^T$$

Error statistics: $E\{\epsilon\} = 0$ and $E\{\epsilon^2\} = V$

No error distribution.

“Bias” or “Intercept” Term

$$y = m(\mathbf{x}; \mathbf{w}) + \epsilon$$

$$m(\mathbf{x}; \mathbf{w}) = \mathbf{b}(\mathbf{x})^T \mathbf{w}, \text{ for } \mathbf{b}(\mathbf{x}) = [1, \phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})]^T$$

$$m(\mathbf{x}; \mathbf{w}) = \sum_{m=1}^M \phi_m(\mathbf{x}) w_{m+1} + w_1$$

$$\text{Train-MSE}(\mathbf{w}) = \frac{1}{N_{\text{TR}}} \sum_{n=1}^{N_{\text{TR}}} \left| y_{\text{TR}}(n) - \left(\sum_{m=1}^M \phi_m(\mathbf{x}) w_{m+1} + w_1 \right) \right|^2$$

Stationarity condition: Train-MSE is minimized at \mathbf{w} where the gradient is 0.

Bias term

This means, that all partial derivatives are 0

$$\Rightarrow \frac{\partial}{\partial w_1} \frac{1}{N_{TR}} \sum_{n=1}^{N_{TR}} \left| y_{TR}(n) - \sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} + w_1 \right|^2 = 0$$

Bias term

$$\begin{aligned} 0 &= \frac{\partial}{\partial w_1} \frac{1}{N_{TR}} \sum_{n=1}^{N_{TR}} \left(\left(y_{TR}(n) - \sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right)^2 + w_1^2 + 2 \left(y_{TR}(n) - \sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right) w_1 \right) \\ &= \frac{1}{N_{TR}} \sum_{n=1}^{N_{TR}} \left(\frac{\partial}{\partial w_1} \left(y_{TR}(n) - \sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right)^2 + \frac{\partial}{\partial w_1} w_1^2 + 2 \frac{\partial}{\partial w_1} \left(y_{TR}(n) - \sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right) w_1 \right) \\ &= \frac{1}{N_{TR}} \sum_{n=1}^{N_{TR}} \left(2w_1 + 2 \left(y_{TR}(n) - \sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right) \right) = 2w_1 + \frac{2}{N_{TR}} \sum_{n=1}^{N_{TR}} \left(y_{TR}(n) - \sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right) \end{aligned}$$

Bias term

$$\Rightarrow 2w_1 + \frac{2}{N_{TR}} \sum_{n=1}^{N_{TR}} \left(y_{TR}(n) - \sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right) = 0$$

$$\Rightarrow w_1 = \frac{1}{N_{TR}} \sum_{n=1}^{N_{TR}} \left(\sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right) - \frac{1}{N_{TR}} \sum_{n=1}^{N_{TR}} y_{TR}(n)$$

$$\Rightarrow w_1 = \frac{1}{N_{TR}} \sum_{n=1}^{N_{TR}} \left(\sum_{m=1}^M \phi_m(\mathbf{x}_{TR}(n)) w_{m+1} \right) - \frac{1}{N_{TR}} \sum_{n=1}^{N_{TR}} y_{TR}(n)$$

“Bias” term w_1 compensates for difference between the average output and the average weighted sum of basis functions. Caution: this is **not the model estimation bias that we saw earlier**.

Maximum Likelihood Parameters

From MSE to LS – Reminder

Let D describe the data model.

Ideally, we'd like to have the solution to

$$\min_{\mathbf{w} \in \mathbb{R}^{M+1}} E_{y,x} \{ |y - m(\mathbf{x}; \mathbf{w})|^2 \}$$

Instead, we estimate m by \hat{m} and the MSE as

$$\min_{\mathbf{w} \in \mathbb{R}^{M+1}} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2$$

where $\mathbf{y} = [y_1, \dots, y_N]$ and $\mathbf{H} = [\mathbf{b}(\mathbf{x}_1), \mathbf{b}(\mathbf{x}_2), \dots, \mathbf{b}(\mathbf{x}_N)]$ depend on training data and \hat{m} .

Least **Squares** (LS) estimates Mean **Squared** Error (MSE)

MSE (cont'd)

What about

$$\min_{\mathbf{w} \in \mathbb{R}^{M+1}} E_{y,x} \{ |y - m(\mathbf{x}; \mathbf{w})| \}$$

It could be estimated by

$$\min_{\mathbf{w} \in \mathbb{R}^{M+1}} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_1$$

In general, why model dissimilarity as squared difference and not anything else?

Seems arbitrary...

Data Model Revisited

$$y = m(\mathbf{x}; \mathbf{w}) + \epsilon$$

Up until now, we had $E\{\epsilon\} = 0$ and $E\{\epsilon^2\} = V$.

This time, we also assume specific **error distribution**: $\epsilon \sim N(0, V)$

Why Gaussian?

Gaussian Distribution

$$X \sim N(\mu, \sigma^2)$$

$$f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right)$$

Central Limit Theorem (CLT): Under common conditions, the sum of many random variables will have an approximately normal distribution.

The normal distribution being the distribution with maximum entropy for a given mean and variance
That is, it makes the fewest assumptions on the RV.

Distribution of Output | Input

$$y = \mathbf{b}(x)^T \mathbf{w} + \epsilon$$

where $\epsilon \sim N(0, V)$

Then, for any given x and \mathbf{w} , y is distributed as:

$$N(\mathbf{b}(x)^T \mathbf{w}, V)$$

Our goal remains to find \mathbf{w}

We will use **training data** (same as before), but also our **assumption on distribution** of $y|X, \mathbf{w}$.

Likelihood

Likelihood is defined on a measurement, drawn from a distribution.

It measures “how likely this particular measurement is” in view of the distribution

For given model (\hat{m} and \mathbf{w}) and any given input (\mathbf{x}_i), how likely is a particular output (y_i)?

$$f_*(y_i|\mathbf{x}_i, \mathbf{w})$$

Conditional PDF of y_i **given** \mathbf{x}_i and \mathbf{w} or the PDF of the conditioning defined RV $y_i|\mathbf{x}_i, \mathbf{w}$

Maximum Likelihood Estimation

MLE Approach: Estimate $\hat{\mathbf{w}}$ so that training outputs $\mathbf{y} = [y_1, y_2, \dots, y_{N_{TR}}]^T$ appear to be the most likely ones, for the given training inputs $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_{TR}}]$.

We know that $(y_i | \mathbf{x}_i, \mathbf{w}) \sim N(\mathbf{w}^T \mathbf{b}(\mathbf{x}_i), V)$. That is:

$$f_*(y_i | \mathbf{x}_i, \mathbf{w}) = \frac{1}{\sqrt{2\pi V}} \exp\left(-\frac{1}{2} \frac{(y_i - \mathbf{b}(\mathbf{x}_i)^T \mathbf{w})^2}{V}\right)$$

But what about the PDF $\mathbf{y} | \mathbf{X}, \mathbf{w} \sim$?

Maximum Likelihood Estimation

We assume that, given the inputs, the outputs are statistically independent. That is,

$$f_*(\mathbf{y}|\mathbf{X}, \mathbf{w}) = f_*(\mathbf{y}|\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_{TR}}\}, \mathbf{w}) = \prod_{i=1}^{N_{TR}} f_*(y_i|\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_{TR}}\}, \mathbf{w})$$

We also assume that y_i is independent of \mathbf{x}_j for $i \neq j$. That is,

$$f_*(y_i|\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_{TR}}\}, \mathbf{w}) = f_*(y_i|\mathbf{x}_i, \mathbf{w})$$

Overall:

$$f_*(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{i=1}^{N_{TR}} f_*(y_i|\mathbf{x}_i, \mathbf{w})$$

Maximum Likelihood Estimation

$$f_*(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{i=1}^{N_{TR}} \frac{1}{\sqrt{2\pi V}} \exp\left(-\frac{1}{2} \frac{(y_i - \mathbf{b}(\mathbf{x}_i)^T \mathbf{w})^2}{V}\right) = \frac{1}{(\sqrt{2\pi V})^{N_{TR}}} \exp\left(-\frac{1}{2} \sum_{i=1}^{N_{TR}} \frac{(y_i - \mathbf{b}(\mathbf{x}_i)^T \mathbf{w})^2}{V}\right)$$

$$f_*(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \frac{1}{(\sqrt{2\pi V})^{N_{TR}}} \exp\left(-\frac{1}{2V} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2\right) = N(\mathbf{H}^T \mathbf{w}, V \mathbf{I}_N)$$

Note: $\mathbf{x} \sim N(\mathbf{m}, \mathbf{C}) \Leftrightarrow f(\mathbf{x}) = \det(2\pi \mathbf{C})^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})\right)$

Maximum Likelihood Estimation

MLE: Estimate $\hat{\mathbf{w}}$ so that training outputs $\mathbf{y} = [y_1, y_2, \dots, y_{N_{TR}}]^T$ appear to be the most likely ones, given training inputs $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N_{TR}}]$.

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w} \in \mathbb{R}^{M+1}} f_{\text{PDF}}(\mathbf{y} | \mathbf{X}, \mathbf{w})$$

So, we want to solve:

$$\max_{\mathbf{w} \in \mathbb{R}^{M+1}} \frac{1}{(\sqrt{2\pi V})^{N_{TR}}} \exp \left(-\frac{1}{2V} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2 \right)$$

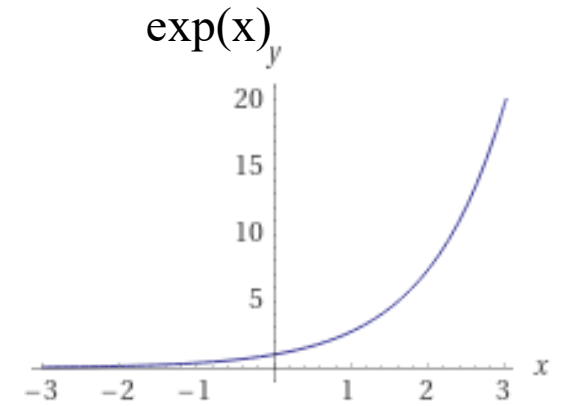
Note on Optimization

- **Two metrics are minimized/maximized by the same argument**, if one is a monotonically increasing function of the other.
 - Examples of monotonically increasing functions: $h(x) = 10C(x)$, $h(x) = C(x)^2$, $h(x) = \sqrt{C(x)}$, $h(x) = \log(C(x))$, $h(x) = \exp(C(x))$.
 - In cases like the above, $\min_w C(x)$ and $\min_w h(x)$ are called equivalent problems.
 - When $h(x)$ is a decreasing function of $C(x)$, then the arguments that minimize $h(x)$ also maximize $C(x)$ (and vice versa).
 - Examples of decreasing functions: $h(x) = -10C(x)$, $h(x) = C(x)^{-2}$, etc.
-

Maximum Likelihood Estimation

Accordingly, the ML problem

$$\max_{\mathbf{w} \in \mathbb{R}^{M+1}} \frac{1}{(\sqrt{2\pi V})^{N_{TR}}} \exp\left(-\frac{1}{2V} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2\right)$$



is equivalent to the simpler:

$$\min_{\mathbf{w} \in \mathbb{R}^{M+1}} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2$$

Which means:

$$\text{MLE} \stackrel{\text{Gaussian Error}}{\equiv} \text{LS}$$

MLE

Which means:

$$\text{MLE}^{\text{Gaussian Error}} \equiv \text{LS}$$

Which means that MSE-LS has likelihood optimality.

The MSE-LS estimate is the most likely one.

For the reasonable case of Normal noise, MSE-LS is **maximum-likelihood optimal**.

What about other distribution assumptions on ϵ ? What if Gaussian but correlated across training data?

Shortcomings of MLE

- Same as MSE – LS
 - To control overfitting, you must increase N or decrease M
-

Bayesian Regression

Bayes Rule

For events A and B:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Offers a way to swap conditioning.

For discrete/continuous random variables $\underline{X} \in S_X, \underline{Z} \in S_Z$, simplify PMF/PDF notation as

$$p_{\underline{X}}(X) \text{ or } f_{\underline{X}}(X) \rightarrow f_*(X)$$

$$p_{\underline{X}, \underline{Z}}(X, Z) \text{ or } f_{\underline{X}, \underline{Z}}(X, Z) \rightarrow f_*(X, Z)$$

$$p_{\underline{X}|\underline{Z}}(X, Z) \text{ or } f_{\underline{X}|\underline{Z}}(X, Z) \rightarrow f_*(X|Z)$$

Bayes Rule becomes:

$$f_*(X|Z) = \frac{f_*(Z|X)f_*(X)}{f_*(Z)} \quad \forall (X, Z) \in S_X \times S_Z$$

Random Parameters

Until now parameter vector \mathbf{w} was seen as deterministic optimization argument.

The best configuration of \mathbf{w}_{best} is unknown. Therefore, we can see it as a random vector.

Thus, it must have some density function, $f_*(\mathbf{w})$.

If f_* was known, how would you choose your parameters?

Choose \mathbf{w} that exhibits the highest probability to be the best.

That is, choose the argument that solves:

$$\max_{\mathbf{w}} f_*(\mathbf{w})$$

Prior Distribution

Can we assume that f_* (the PDF of the best parameter configuration) is known?

Probably not really.

But we can assume that we have a good initial (=we call it “prior”) guess of it.

From now on, f_* denotes our prior guess on the distribution of the best \mathbf{w} .

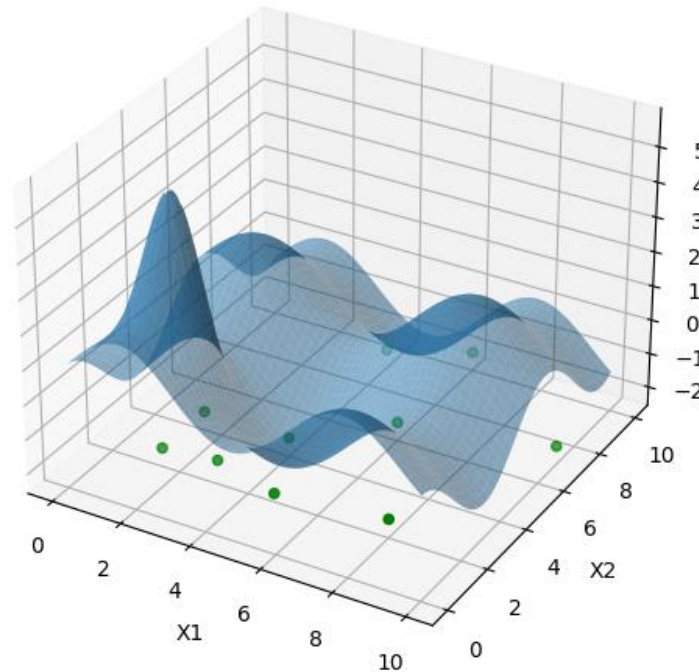
For example, we can assume

$$f_*(\mathbf{w}) = N(\mathbf{m}_0, \mathbf{\Sigma}_0)$$

for some $\mathbf{m}_0, \mathbf{\Sigma}_0$, since Gaussian makes the fewest assumptions.

Even simpler, we can assume zero-mean isotropic Gaussian $f_*(\mathbf{w}) = N(\mathbf{0}, aI)$

Prior Distribution (cont'd)



Example: Considering GRBF model, if we have a guess or expert knowledge that in our data model output y takes higher values when the input is around $[2, 2]^T$, then I can assume that the parameters weighting the GRBF-centers that are near $[2, 2]^T$ would take higher values with higher probability.

Prior Distribution (cont'd)

So f_* denotes our prior guess on the distribution of the best \mathbf{w} .

Again, we can choose our parameters to be the solution to $\max_{\mathbf{w}} f_*(\mathbf{w})$

If $f_*(\mathbf{w}) = N(\mathbf{m}_0, \Sigma_0)$, $\hat{\mathbf{w}} = \mathbf{m}_0$, since Gaussian bell attains max value at the mean.

But maybe I can use some training data to help us improve my best-parameter distribution guess.

Let (\mathbf{y}, \mathbf{X}) be the training dataset ($y_i = [\mathbf{y}]_i$ is the output corresponding to input $[\mathbf{X}]_{:,i}$).

Our updated guess of the best-parameter distribution, after training data (\mathbf{y}, \mathbf{X}) have been studied, we call it “posterior” or “a posteriori” distribution.

- **Prior** guess on best-parameter distribution: $f_*(\mathbf{w})$
 - **Posterior** guess on best-parameter distribution: $f_*(\mathbf{w}|\mathbf{y}, \mathbf{X})$
-

Posterior Distribution

But maybe I can use some training data to help us improve my best-parameter distribution guess.

Let (\mathbf{y}, \mathbf{X}) be the training dataset ($y_i = [\mathbf{y}]_i$ is the output corresponding to input $[\mathbf{X}]_{:,i}$).

Our updated guess of the best-parameter distribution, after training data (\mathbf{y}, \mathbf{X}) have been studied, we call it “posterior” or “a posteriori” distribution.

- **Prior** guess on best-parameter distribution: $f_*(\mathbf{w})$
 - **Posterior** guess on best-parameter distribution: $f_*(\mathbf{w}|\mathbf{y}, \mathbf{X})$
-

Maximum A Posteriori Probability

When we had only the prior distribution, we designed our parameters by solving:

$$\max_w f_*(\mathbf{w})$$

If we use the data in (\mathbf{y}, \mathbf{X}) a smart way and derive the posterior distribution/probability $f_*(\mathbf{w}|\mathbf{y}, \mathbf{X})$, then we'll design our parameters by solving:

$$\max_w f_*(\mathbf{w}|\mathbf{y}, \mathbf{X})$$

This is called the **Maximum A Posteriori Probability** approach.

Two steps follow:

Step 1: $f_*(\mathbf{w}|\mathbf{y}, \mathbf{X})$

Step 2: Solve $\max_w f_*(\mathbf{w}|\mathbf{y}, \mathbf{X})$

Posterior Distribution

We are looking for posterior distribution $f_*(\mathbf{w} | \mathbf{y}, \mathbf{X})$.

This means that \mathbf{w} is the variable and that \mathbf{y}, \mathbf{X} are given (conditioned over) and fixed.

Therefore, terms like $f_*(\mathbf{X})$ and $f_*(\mathbf{y}, \mathbf{X})$ are just positive constants with respect to \mathbf{w} .

$$\begin{aligned} f_*(\mathbf{w} | \mathbf{y}, \mathbf{X}) &= \frac{f_*(\mathbf{y}, \mathbf{X} | \mathbf{w}) f_*(\mathbf{w})}{f_*(\mathbf{y}, \mathbf{X})} \propto f_*(\mathbf{y}, \mathbf{X} | \mathbf{w}) f_*(\mathbf{w}) \\ &\propto f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) f_*(\mathbf{X} | \mathbf{w}) f_*(\mathbf{w}) \\ &\propto f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) f_*(\mathbf{X}) f_*(\mathbf{w}) \\ &\propto f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) f_*(\mathbf{w}) \end{aligned}$$

Thus, $f_*(\mathbf{w} | \mathbf{y}, \mathbf{X}) = f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) f_*(\mathbf{w}) \cdot \text{pos_const}$

Posterior Distribution

$$f_*(\mathbf{w} | \mathbf{y}, \mathbf{X}) = f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) f_*(\mathbf{w}) \cdot \text{pos_const}$$

During our MLE studies we we found:

$$f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) = N(\mathbf{H}^T \mathbf{w}, V \mathbf{I}_N) = \text{pos_const} \cdot \exp \left(-\frac{1}{2} (\mathbf{y} - \mathbf{H}^T \mathbf{w})^T (V \mathbf{I}_N)^{-1} (\mathbf{y} - \mathbf{H}^T \mathbf{w}) \right)$$

About our prior, we assume Gaussian (fewest assumptions) with specific mean and covariance:

$$f_*(\mathbf{w}) = N(\mathbf{m}_0, \mathbf{\Sigma}_0) = \text{pos_const} \cdot \exp \left(-\frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T \mathbf{\Sigma}_0^{-1} (\mathbf{w} - \mathbf{m}_0) \right)$$

Posterior Distribution (cont'd)

Thus,

$$f_*(\mathbf{w} | \mathbf{y}, \mathbf{X}) = f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) f_*(\mathbf{w}) \cdot \text{pos_const}$$

$$= \text{pos_const} \cdot \exp \left(\left(-\frac{1}{2} (\mathbf{y} - \mathbf{H}^T \mathbf{w})^T (V \mathbf{I}_N)^{-1} (\mathbf{y} - \mathbf{H}^T \mathbf{w}) \right) + \left(-\frac{1}{2} (\mathbf{w} - \mathbf{m}_w)^T (\mathbf{S}_w)^{-1} (\mathbf{w} - \mathbf{m}_w) \right) \right)$$

$$= \text{pos_const} \cdot \exp \left(-\frac{1}{2} \left((\mathbf{y} - \mathbf{H}^T \mathbf{w})^T (V \mathbf{I}_N)^{-1} (\mathbf{y} - \mathbf{H}^T \mathbf{w}) + (\mathbf{w} - \mathbf{m}_0)^T \mathbf{\Sigma}_0^{-1} (\mathbf{w} - \mathbf{m}_0) \right) \right)$$

Posterior Distribution (cont'd)

Make sure you know how to derive:

$$\begin{aligned} & (\mathbf{y} - \mathbf{H}^T \mathbf{w})^T (\mathbf{V} \mathbf{I}_N)^{-1} (\mathbf{y} - \mathbf{H}^T \mathbf{w}) + (\mathbf{w} - \mathbf{m}_0)^T \Sigma_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &= (\mathbf{y} - \mathbf{H}^T \mathbf{w})^T \mathbf{V}^{-1} (\mathbf{y} - \mathbf{H}^T \mathbf{w}) + (\mathbf{w} - \mathbf{m}_0)^T \Sigma_0^{-1} (\mathbf{w} - \mathbf{m}_0) \\ &= \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{w}^T \mathbf{H}^T \mathbf{V}^{-1} \mathbf{H} \mathbf{w} - 2 \mathbf{w}^T (\mathbf{H} \mathbf{V}^{-1} \mathbf{y}) + \mathbf{w}^T \Sigma_0^{-1} \mathbf{w} + \mathbf{m}_0^T \Sigma_0^{-1} \mathbf{m}_0 - 2 \mathbf{w}^T \Sigma_0^{-1} \mathbf{m}_0 \\ &= \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{m}_0^T \Sigma_0^{-1} \mathbf{m}_0 + \mathbf{w}^T (\mathbf{H} \mathbf{V}^{-1} \mathbf{H}^T + \Sigma_0^{-1}) \mathbf{w} - 2 \mathbf{w}^T (\mathbf{H} \mathbf{V}^{-1} \mathbf{y} + \Sigma_0^{-1} \mathbf{m}_0) \\ &= \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{m}_0^T \Sigma_0^{-1} \mathbf{m}_0 + \mathbf{w}^T \Sigma^{-1} \mathbf{w} - 2 \mathbf{w}^T \mathbf{z} \\ &= \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{m}_0^T \Sigma_0^{-1} \mathbf{m}_0 + \mathbf{w}^T \Sigma^{-1} \mathbf{w} - 2 \mathbf{w}^T \Sigma^{-1} \Sigma \mathbf{z} \\ &= \mathbf{y}^T \mathbf{V}^{-1} \mathbf{y} + \mathbf{m}_0^T \Sigma_0^{-1} \mathbf{m}_0 + \mathbf{w}^T \Sigma^{-1} \mathbf{w} - 2 \mathbf{w}^T \Sigma^{-1} \mathbf{m} + \mathbf{m}^T \mathbf{m} - \mathbf{m}^T \mathbf{m} \\ &= \mathbf{w}^T \Sigma^{-1} \mathbf{w} - 2 \mathbf{w}^T \Sigma^{-1} \mathbf{m} + \mathbf{m}^T \mathbf{m} + \text{pos_const} \\ &= (\mathbf{w} - \mathbf{m})^T \Sigma^{-1} (\mathbf{w} - \mathbf{m}) + \text{pos_const} \end{aligned}$$

Posterior Distribution (cont'd)

$$f_*(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \beta \cdot \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \mathbf{m})\right) = N(\mathbf{m}, \boldsymbol{\Sigma})$$

$\mathbf{w} | \mathbf{y}, \mathbf{X}$ is multivariate Gaussian with mean \mathbf{m} and covariance matrix $\boldsymbol{\Sigma}$, where

$$\boldsymbol{\Sigma} = \left(\frac{1}{V} \mathbf{H} \mathbf{H}^T + \boldsymbol{\Sigma}_0^{-1}\right)^{-1}$$

$$\mathbf{m} = \boldsymbol{\Sigma} \left(\frac{1}{V} \mathbf{H} \quad \mathbf{y} + \boldsymbol{\Sigma}_0^{-1} \mathbf{m}_0\right) = (\mathbf{H} \mathbf{H}^T + V \boldsymbol{\Sigma}_0^{-1})^{-1} (\mathbf{H} \mathbf{y} + \boldsymbol{\Sigma}_0^{-1} \mathbf{m}_0 V)$$

Maximize Posterior Density

$$f_*(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \beta \cdot \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{m})^T \mathbf{\Sigma}^{-1}(\mathbf{w} - \mathbf{m})\right) = N(\mathbf{m}, \mathbf{\Sigma})$$

Gaussian PDF is maximized at its mean.

That is, the MAP parameter vector is

$$\hat{\mathbf{w}} = (\mathbf{H}\mathbf{H}^T + \mathbf{V}\mathbf{\Sigma}_0^{-1})^{-1}(\mathbf{H}\mathbf{y} + \mathbf{\Sigma}_0^{-1}\mathbf{m}_0 \mathbf{V})$$

Read is data, **blue** is prior distribution.

Alternative Derivation

We showed that $\mathbf{w}|\mathbf{y}, \mathbf{X}$ follows $N(\mathbf{m}, \Sigma)$. We recognized that MAP parameters should be \mathbf{m} .

Here is an alternative derivation, without finding first the distribution of $\mathbf{w}|\mathbf{y}, \mathbf{X}$.

$$\begin{aligned} & \max_{\mathbf{w}} f_*(\mathbf{w} | \mathbf{y}, \mathbf{X}) \\ & \equiv \max_{\mathbf{w}} f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) f_*(\mathbf{w}) \cdot \text{pos_onst.} \\ & \equiv \max_{\mathbf{w}} f_*(\mathbf{y} | \mathbf{X}, \mathbf{w}) f_*(\mathbf{w}) \\ & \equiv \max_{\mathbf{w}} \ln(f_*(\mathbf{y}, \mathbf{X} | \mathbf{w})) + \ln(f_*(\mathbf{w})) \end{aligned}$$

Posterior Density Maximization

$$\max_{\mathbf{w}} f_*(\mathbf{w} | \mathbf{y}, \mathbf{X})$$

$$\equiv \max_{\mathbf{w}} \ln \left(\frac{1}{(\sqrt{2\pi V})^N} \exp \left(-\frac{1}{2V} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2 \right) \right) + \ln \left(\frac{1}{\sqrt{2\pi |\Sigma_0|}} \exp \left(-\frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T \Sigma_0^{-1} (\mathbf{w} - \mathbf{m}_0) \right) \right)$$

$$\equiv \max_{\mathbf{w}} \ln \left(\frac{1}{(\sqrt{2\pi V})^N} \right) - \frac{1}{2V} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2 + \ln \left(\frac{1}{\sqrt{2\pi |\Sigma_0|}} \right) - \frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T \Sigma_0^{-1} (\mathbf{w} - \mathbf{m}_0)$$

$$\max_{\mathbf{w}} f_*(\mathbf{w} | \mathbf{y}, \mathbf{X}) \equiv \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2 + (\mathbf{w} - \mathbf{m}_0)^T (V \Sigma_0^{-1}) (\mathbf{w} - \mathbf{m}_0)$$

Posterior Density Maximization (cont'd)

$$\max_{\mathbf{w}} \|\mathbf{y} - \mathbf{H}^T \mathbf{w}\|_2^2 + (\mathbf{w} - \mathbf{m}_0)^T (V \Sigma_0^{-1}) (\mathbf{w} - \mathbf{m}_0)$$

Convex objective function. Apply stationarity condition:

$$C(\mathbf{w}) = \mathbf{y}^T \mathbf{y} + \mathbf{w}^T \mathbf{H} \mathbf{H}^T \mathbf{w} - 2 \mathbf{w}^T \mathbf{H} \mathbf{y} + \mathbf{w}^T (V \Sigma_0^{-1}) \mathbf{w} + \mathbf{m}_0^T (V \Sigma_0^{-1}) \mathbf{m}_0 - 2 \mathbf{w}^T (V \Sigma_0^{-1}) \mathbf{m}_0$$

$$\mathbf{g}(\mathbf{w}) = 2 \mathbf{H} \mathbf{H}^T \mathbf{w} - 2 \mathbf{H} \mathbf{y} + 2 (V \Sigma_0^{-1}) \mathbf{w} - 2 (V \Sigma_0^{-1}) \mathbf{m}_0$$

$$\mathbf{g}(\hat{\mathbf{w}}) = \mathbf{0} \Leftrightarrow$$

$$\hat{\mathbf{w}} = (\mathbf{H} \mathbf{H}^T + V \Sigma_0^{-1})^{-1} (\mathbf{H} \mathbf{y} + \Sigma_0^{-1} \mathbf{m}_0 V)$$

MAP – Special Cases

MAP solution:

$$\hat{\mathbf{w}} = (\mathbf{H}\mathbf{H}^T + \mathbf{V}\boldsymbol{\Sigma}_0^{-1})^{-1}(\mathbf{H}\mathbf{y} + \boldsymbol{\Sigma}_0^{-1}\mathbf{m}_0 \mathbf{V})$$

No data? Boils down to prior distribution maximization.

$$\mathbf{H} = \mathbf{0} \Rightarrow \hat{\mathbf{w}} = (\mathbf{0} + \mathbf{V}\boldsymbol{\Sigma}_0^{-1})^{-1}(\mathbf{0} + \boldsymbol{\Sigma}_0^{-1}\mathbf{m}_0 \mathbf{V}) = \mathbf{m}_0$$

No prior distribution assumption? Boils down to MLE/MSE/LS.

Equivalent to considering $\boldsymbol{\Sigma}_0 = a\mathbf{I}$ with $a \rightarrow \infty$.

$$a \rightarrow \infty \Rightarrow \boldsymbol{\Sigma}_0^{-1} = \mathbf{0} \Rightarrow \hat{\mathbf{w}} = (\mathbf{H}\mathbf{H}^T + \mathbf{0})^{-1}(\mathbf{H}\mathbf{y} + \mathbf{0}) = (\mathbf{H}\mathbf{H}^T)^{-1}\mathbf{H}\mathbf{y}$$

MLE vs MAP

- Same as comparison between LS and Regularized-LS.
 - MLE relies solely on data and it needs many, otherwise might overfit.
 - MAP with Gaussian error and prior is LS + regularization.
 - MAP relies also on model priors; less prone to overfit; could work with less data.
 - Prior with very high variance: MAP tends to MLE/MSE/LS
 - Limited/no data? MAP tends to prior mean.
-