

BAYESIAN LEARNING

- There is a natural interpretation of the regularization procedure in terms of the Bayesian approach:

Starting point, same as before, the likelihood:

$$p(D|\vec{\theta}) = \prod_{i=1}^N \mathcal{N}(y_i | f(x_i, \vec{\theta}), \sigma^2)$$

But now we need a prior distribution for $\vec{\theta}$:

Suppose $p(\vec{\theta}) = \mathcal{N}(\vec{\theta} | \vec{0}, \alpha^{-1} \mathbf{1})$

$$= \left(\frac{\alpha}{2\pi}\right)^M \exp\left\{-\frac{\alpha}{2} \vec{\theta}^T \cdot \vec{\theta}\right\}$$

Then the posterior

$$M = \dim \vec{\theta}$$

$$p(\vec{\theta}|D) \propto p(D|\vec{\theta}) p(\vec{\theta})$$

$$\ln p(\vec{\theta}|D) = \ln p(D|\vec{\theta}) + \ln p(\vec{\theta}) + \text{const.}$$

$$-\ln p(\vec{\theta} | D) = \frac{1}{2\sigma^2} \sum_i (y_i - f(x_i, \vec{\theta}))^2 + \frac{\alpha}{2} \vec{\theta}^T \cdot \vec{\theta} + \text{const}$$

$$= \frac{1}{\sigma^2} \left\{ \frac{1}{2} \sum_i (y_i - f(x_i, \vec{\theta}))^2 + \frac{\tilde{\alpha}}{2} \vec{\theta}^T \cdot \vec{\theta} \right\}$$

where

$$\tilde{\alpha} \equiv \alpha \sigma^2$$

= E_{ridge} !

So we see that obtaining the $\vec{\theta}_{\text{MAP}}$ in the Bayesian approach is equivalent to $\vec{\theta}_{\text{ridge}}$, for a Gaussian choice of the prior.

Bayesian Linear Regression

We have just seen that the Bayesian approach automatically encodes a regularized fitting. However its most attractive feature is that it provides a natural way to estimate prediction uncertainties (through the "predictive distribution" we saw in lecture #2)

$$P(y_{\text{new}} | \vec{x}_{\text{new}}, D) = \int d\vec{\theta} p(\vec{\theta} | D) p(y_{\text{new}} | \vec{x}_{\text{new}}, \vec{\theta})$$

Remember that in order to compute $p(\vec{\theta} | D)$ in principle you need to compute the normalisation (the "evidence"):

$$\int d\vec{\theta} p(D | \vec{\theta}) p(\vec{\theta})$$

There are only few cases where these integrals are analytical, and in real-life problems this is typically not the case. The most popular academic example is:

Likelihood $\left\{ \begin{array}{l} p(y | \vec{x}, \vec{\theta}) = \mathcal{N}(y | f(\vec{x}, \vec{\theta}), \Sigma) \\ \text{1 point} \end{array} \right.$ (assume known)

$$f(\vec{x}, \vec{\theta}) = \vec{\theta}^T \cdot \vec{\phi}(\vec{x}) = \theta_1 \phi_1(\vec{x}) + \dots + \theta_M \phi_M(\vec{x})$$

• Prior: $p(\vec{\theta}) = \mathcal{N}(\vec{\theta} | \vec{\mu}_0, \Sigma_0)$

So the posterior $p(\vec{\theta} | \vec{y})$ and the 'evidence' $p(\vec{y})$ are both Gaussians,

only because since $f(x; \theta)$ is linear in θ , then in the exponent we will have a polynomial of order 2 in θ , which can be recasted as a Gaussian by completing the square. So the result is

$$p(\vec{\theta} | \vec{y}) = \mathcal{N}(\vec{\theta} | \vec{\mu}_N, S_N)$$

$$\text{where } \vec{\mu}_N = S_N \{ \Phi^T \bar{C}^{-1} \vec{y} + S_0^{-1} \vec{\mu}_0 \}$$

$$S_N = (\Phi^T \bar{C}^{-1} \Phi + S_0^{-1})^{-1}$$

for the posterior, while for the evidence is:

$$p(\vec{y}) = \mathcal{N}(\vec{y} | \Phi \vec{\mu}_0, \bar{C} + \Phi S_0 \Phi^T)$$

$N \times 1$

$(N \times D) \times (D \times 1)$

$N \times N$

$(N \times D) \times (D \times D) \times (D \times N)$

In an analogous fashion to the evidence, the predictive distribution is an integral over two Gaussians, and it's found to be:

$$p(y_{\text{new}} | \vec{x}_{\text{new}}, \mathcal{D}) = \mathcal{N}(y_{\text{new}} | \vec{\mu}_N^T \cdot \vec{\phi}(\vec{x}_{\text{new}}), \sigma_N^2(\vec{x}_{\text{new}}))$$

where $\sigma_N^2(\vec{x}_{\text{new}}) = \sigma^2(x_{\text{new}}) + \vec{\phi}(\vec{x}_{\text{new}})^T S_N \vec{\phi}(\vec{x}_{\text{new}})$

statistical noise at this point

So even if σ^2 (the data noise is the same $\forall \vec{x}$, still σ_N^2 would get a dependence on \vec{x} coming from the 2nd term. As we will see, this dependence is such that the variance grows for points \vec{x}_{new} farther away from any of the \vec{x}_i considered in the dataset, and viceversa: the variance of the predictive

distribution shrinks as \vec{x}_{new} gets closer to the \vec{x}_i . This is intuitive: as we get farther away from the points used in the fit, our predictions become more uncertain.

[Show notebook on Bayesian prediction]

- Before understanding analytically this behaviour with the # of observations (N points), let us visualise in another way the sequential nature of the Bayesian learning:

The posterior distribution after seeing N datapoints

$$p(\vec{w} | \vec{y}^{(N)}) \propto \underbrace{p(\vec{y}^{(N)} | \vec{w})}_{\text{joint likelihood of } N \text{ observations}} p(\vec{w}) \quad \left. \vphantom{p(\vec{w} | \vec{y}^{(N)})} \right\} \text{prior before any observation}$$

$$\vec{y}^{(N)} = \{y_1, y_2, \dots, y_N\}$$

$$\begin{aligned}
 &= \underbrace{p(y_N | \vec{w})}_{\text{likelihood of the } N^{\text{th}} \text{ observation}} \underbrace{p(\vec{y}_{(N-1)} | \vec{w})}_{\text{joint likelihood of the rest of } N-1 \text{ observations}} p(\vec{w}) \\
 &\propto p(y_N | \vec{w}) \underbrace{p(\vec{w} | \vec{y}_{(N-1)})}_{\text{posterior after seeing } N-1 \text{ observations}} \stackrel{\text{prior before the } N^{\text{th}} \text{ observation}}{=}
 \end{aligned}$$

So the posterior after one observation becomes the prior to the 2nd observation, and so on...

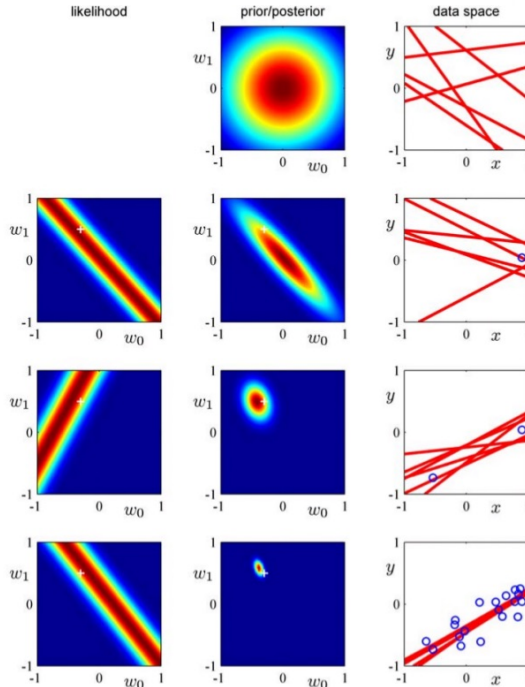


Figure 3.7 Illustration of sequential Bayesian learning for a simple linear model of the form $y(x, w) = w_0 + w_1 x$. A detailed description of this figure is given in the text.

- How does the posterior of \vec{w} behaves with the # of points, N ?

$$p(\vec{\theta} | \vec{y}) = \mathcal{N}(\vec{\theta} | \vec{\mu}_N, S_N)$$

just for simplicity.
 → Results in general still valid

Assume covariance matrix of likelihood $C = \sigma^2 \mathbf{1}$

$$\vec{\mu}_N = S_N \{ \sigma^{-2} \Phi^T \vec{y} + S_0^{-1} \vec{\mu}_0 \}$$

$$S_N = (\sigma^{-2} \Phi^T \Phi + S_0^{-1})^{-1}$$

- 1) $N=0$ (i.e. no observations yet)

$$S_N = (0 + S_0^{-1})^{-1} = S_0$$

$$\vec{\mu}_N = S_N \{ 0 + S_0^{-1} \vec{\mu}_0 \} = \vec{\mu}_0$$

↳ So you get back the prior

- 2) $N \rightarrow \infty$

$\Phi^T \Phi$ grows with N

{ This matrix is less and less singular

also $\Phi^T \vec{y}$ grows

{ more terms in the scalar product, which grows modulo fine tuning

So $S_N \approx (\sigma^{-2} \Phi^T \Phi)^{-1} \Rightarrow$ variance decreases with N

$$\vec{\mu}_N \approx S_N \{ \sigma^{-2} \Phi^T \vec{y} \}$$

$$\approx (\Phi^T \Phi)^{-1} \Phi^T \vec{y} \quad \text{as } N \rightarrow \infty$$

\hookrightarrow The mean of the posterior tends to the MLE solution! \rightarrow No prior dependence

• Bayesian Model Comparison

Even though the introduction of a prior plays the role of a regularizer when doing the fit, and a highly parametrised model can be fitted to some simple data without doing overfit, this may still not be a good model:

In Bayesian approach we use the Evidence (denominator of the Bayes theorem) to compare among different models

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} \rightarrow p(D|\text{model})$$

So imagine you have 2 models m_1 & m_2 ,
with Evidences $p(D|m_1)$ and $p(D|m_2)$

$$p(m_i|D) \propto p(D|m_i) p(m_i)$$

← prior of the model itself
(not of its parameters!)

If all models are assumed

to have the same prior, then the most probable
model among m_1 and m_2 is obtained by
computing the "Bayes factor"

$$B \equiv p(D|m_1) / p(D|m_2)$$

for a sufficiently large B (e.g. > 10) one says
there is a strong "evidence" in favour of m_1

Finally, let's get an intuition about which are the features of a model which may penalize it from the evidence point of view.

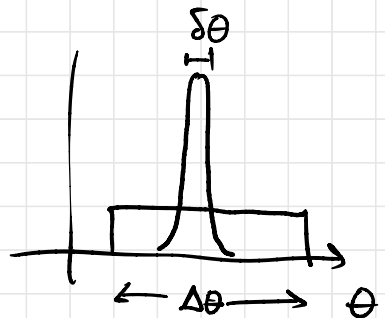
$$p(D|m) = \int d\vec{\theta} p(D|\vec{\theta}, m) p(\vec{\theta}|m)$$

Suppose the following simplifications (which however do not change the general picture)

- Only 1 parameter
- likelihood sharply peaked around $\theta = \theta^*$, with width = $\delta\theta$
- flat prior $p(\theta) = \frac{1}{\Delta\theta}$

$$p(D|m) \approx \frac{\delta\theta}{\Delta\theta} p(D|\theta^*)$$

* Penalization of models having large prior uncertainties



• Now assume we have P parameters

$$\vec{\theta} = \{\theta_k\}_{k=1}^P$$

Suppose a similar ratio

$$\frac{\delta\theta_k}{\Delta\theta_k} \quad \forall k$$

Then the evidence

$$p(D|m) \approx p(D|\vec{\theta}^*) \left[\frac{\delta\theta}{\Delta\theta} \right]^M$$

✦ Complex models with many parameters are penalized because of the 2nd factor. However typically these larger models fit better the data, so $p(D|\vec{\theta}^*)$ is larger.

⇒ Optimal model should have an intermediate complexity