# Bayesian model comparison

Model selection from a Bayesian perspective: use probability theory to represent uncertainties in the model choice.

Namely, recall that

$$\overbrace{p(\vec{\omega}|D, M_i)}^{\text{posterior}} = \frac{p(D|\vec{\omega}, M_i)\, p(\vec{\omega}|M_i)}{p(D|M_i)}, \quad \text{where}$$

↑ data

model out of a set of $L$ models: $\{M_i\}$ that we wish to compare

$$p(D|M_i) = \int d\vec{\omega}\, p(D|\vec{\omega}, M_i)\, p(\vec{\omega}|M_i)$$

is the model evidence.

Now, we can consider

$$\underbrace{p(M_i|D)}_{\text{model posterior}} \sim p(D|M_i)\, \underbrace{p(M_i)}_{\substack{\text{often, equal}\\ \text{probs for each}\\ \text{model}}}$$

$\dfrac{p(D|M_i)}{p(D|M_j)}$ is called a Bayes factor

Finally, the predictive distr'n is given by

$$p(t \mid \vec{x}, D) = \sum_{i=1}^{L} \overbrace{p(t \mid \vec{x}, M_i, D)}^{\text{predictive distr'n of model } M_i} p(M_i \mid D) \qquad (\ast\ast)$$
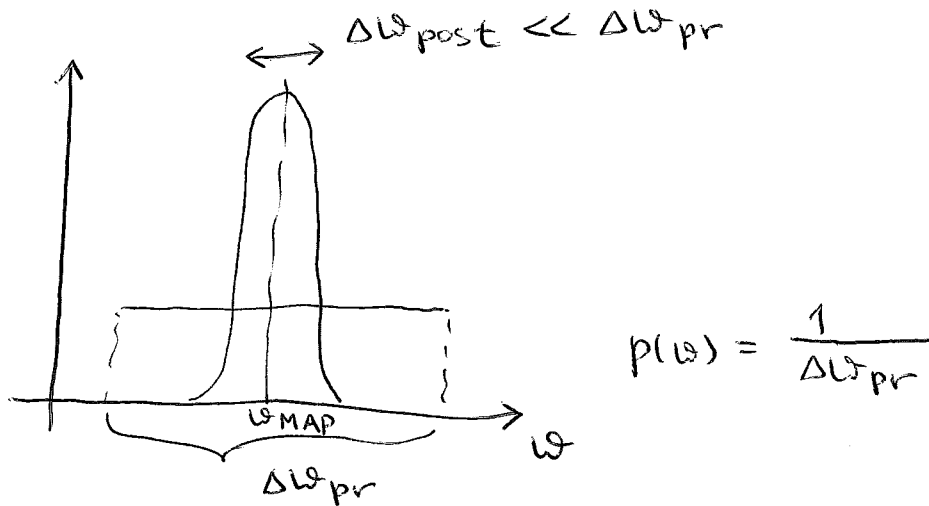
a mixture distr'n weighted by $p(M_i \mid D)$

a simple appr'n to $(\ast\ast)$ is to leave only a single term in the sum corresponding to the most probable model for which $p(M_i \mid D)$ is at maximum. $\leftarrow$ (model selection)

---

Consider a model $M_i$ with a single prm $\omega$.

Then $p(\omega \mid D) \sim p(D \mid \omega) p(\omega)$ [dependence on $M_i$ omitted for brevity]

Assume that $p(\omega \mid D)$ is peaked @ $\omega_{MAP}$ with width $\Delta\omega_{post}$ and the prior is flat with width $\Delta\omega_{pr}$:



$$p(\omega) = \frac{1}{\Delta\omega_{pr}}$$

Then $p(D) = \underbrace{\int d\omega \, p(D \mid \omega) p(\omega)}_{\substack{\text{evidence} \\ \text{for } M_i}} \simeq p(D \mid \omega_{MAP}) \dfrac{\Delta\omega_{post}}{\Delta\omega_{pr}},$

or

$$\log P(D) \simeq \log P(D \mid \omega_{MAP}) + \log\left(\frac{\Delta\omega_{post}}{\Delta\omega_{pr}}\right)$$

$< 0$, acts like a penalty for model complexity

For a model with $M$ prms,

$$\log P(D) \simeq \log P(D \mid \omega_{MAP}) + M \log\left(\frac{\Delta\omega_{post}}{\Delta\omega_{pr}}\right)$$

if $\frac{\Delta\omega_{post}}{\Delta\omega_{pr}}$ is the same for each prm.

So, the penalty $\sim M$.

---

Next, consider 3 models $M_1, M_2, M_3$ of increasing complexity. Imagine generating synthetic data sets from each of these models: choose prms from $p(\vec{\omega})$ & then generate data from $p(D \mid \vec{\omega}) \Rightarrow$ compute $p(D)$. (or, more precisely, $p(D \mid M_i)$, $i = 1, 2, 3$)



$p(D)$

$M_1$ is least flexible, produces "similar" datasets

$M_2$ is intermediate

$M_3$ is most flexible

"1D projection of dataset $D$"

$D_o$
observed dataset

$\begin{cases} p(D \mid M_2) > p(D \mid M_1) \\ p(D \mid M_2) > p(D \mid M_3) \end{cases}$   $M_1$ did not fit well, too inflexible

$M_3$ is too flexible

$\boxed{M_2 \text{ wins}}$

# The evidence approximation

Consider

$$p(t|\vec{t}) = \int d\vec{w}\, d\alpha\, d\beta\, \underbrace{p(t|\vec{w},\beta)}_{\text{likelihood}} \overbrace{p(\vec{w}|\vec{t},\alpha,\beta)}^{\text{posterior for }\vec{w}} \overbrace{p(\alpha,\beta|\vec{t})}^{\text{posterior for }\alpha,\beta}$$

$$\underbrace{\underbrace{p(t,\vec{w},\alpha,\beta|\vec{t})}\qquad\overbrace{p(\vec{w},\alpha,\beta|\vec{t})}}$$

Omitted dependence on $\vec{x}$, $X$ for brevity

If $p(\alpha,\beta|\vec{t})$ is sharply peaked around $\hat{\alpha}, \hat{\beta}$, then

$$p(t|\vec{t}) \simeq p(t|\vec{t},\hat{\alpha},\hat{\beta}) = \int d\vec{w}\, p(t|\vec{w},\hat{\beta})\, p(\vec{w}|\vec{t},\hat{\alpha},\hat{\beta})$$

Further, $p(\alpha,\beta|\vec{t}) \sim p(\vec{t}|\alpha,\beta)\, p(\alpha,\beta)$

If $p(\alpha,\beta)$ is flat, we can simply maximize $p(\vec{t}|\alpha,\beta)$ to find $\hat{\alpha},\hat{\beta}$.

Now, $p(\vec{t}|\alpha,\beta) = \int d\vec{w}\, p(\vec{t}|\vec{w},\beta)\, p(\vec{w}|\alpha) = $

$$= \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}}\left(\frac{\alpha}{2\pi}\right)^{\frac{D+1}{2}} \int d\vec{w}\, e^{-E(\vec{w})}, \quad \text{where}$$

$$E(\vec{w}) = \frac{\beta}{2}\underbrace{\sum_{n=1}^{N}(t_n - \vec{w}^T\vec{\varphi}(\vec{x}_n))^2}_{\|\vec{t} - \Phi\vec{w}\|^2} + \frac{\alpha}{2}\vec{w}^T\vec{w} \quad (\equiv)$$

$$\overbrace{(\vec{t}-\Phi\vec{m}_N)^T(\vec{t}-\Phi\vec{m}_N)}$$

$$\equiv \underbrace{\frac{\beta}{2}\|\vec{t}-\Phi\vec{m}_N\|^2 + \frac{\alpha}{2}\vec{m}_N^T\vec{m}_N}_{E(\vec{m}_N)} +$$

$$+ \frac{1}{2}(\vec{w}-\vec{m}_N)^T\underbrace{(\alpha I + \beta\Phi^T\Phi)}_{A = S_N^{-1}}(\vec{w}-\vec{m}_N).$$

Here, $\vec{m}_N = \beta A^{-1}\Phi^T\vec{t} = \beta S_N\Phi^T\vec{t}$

Indeed, the last term on the RHS is

$$\frac{\alpha}{2}\vec{w}^T\vec{w} - \alpha\vec{w}^T\vec{m}_N + \frac{\alpha}{2}\vec{m}_N^T\vec{m}_N +$$

$$+ \frac{\beta}{2}\vec{w}^T\Phi^T\Phi\vec{w} + \frac{\beta}{2}\vec{m}_N^T\Phi^T\Phi\vec{m}_N -$$

$$- \frac{\beta}{2}\vec{w}^T\Phi^T\Phi\vec{m}_N - \frac{\beta}{2}\vec{m}_N^T\Phi^T\Phi\vec{w} , \text{ so that}$$

$$E(\vec{w}) = \frac{\beta}{2}(\vec{t}-\Phi\vec{m}_N)^T(\vec{t}-\Phi\vec{m}_N) + \frac{\alpha}{2}\vec{w}^T\vec{w} -$$

$$- \alpha\vec{w}^T\vec{m}_N + \alpha\vec{m}_N^T\vec{m}_N + \frac{\beta}{2}\vec{w}^T\Phi^T\Phi\vec{w} +$$

$$+ \frac{\beta}{2}\vec{m}_N^T\Phi^T\Phi\vec{m}_N - \beta(\Phi\vec{w})^T(\Phi\vec{m}_N) =$$

$$= \frac{\beta}{2}\vec{t}^T\vec{t} + \frac{\beta}{2}(\Phi\vec{m}_N)^T(\Phi\vec{m}_N) - \beta(\Phi\vec{m}_N)^T\vec{t} +$$

$$+ \frac{\alpha}{2}\vec{w}^T\vec{w} - \ldots =$$

$$= \frac{\beta}{2}\vec{t}^T\vec{t} + \frac{\beta}{2}(\Phi\vec{w})^T(\Phi\vec{w}) + \frac{\alpha}{2}\vec{w}^T\vec{w}$$

OK, $\theta(\vec{w}^2)$: "new": $\frac{1}{2}\vec{w}^T A \vec{w} =$

terms

$$= \frac{1}{2}\vec{w}^T(\alpha I + \beta \phi^T\phi)\vec{w} =$$

$$= \frac{1}{2}\vec{w}^T \alpha \vec{w} + \frac{\beta}{2}(\phi\vec{w})^T(\phi\vec{w}).$$

"old": $\frac{\beta}{2}(\phi\vec{w})^T(\phi\vec{w}) + \frac{\alpha}{2}\vec{w}^T\vec{w}$    same!

$$\underline{\underline{\quad}}$$

$\theta(\vec{w}^0)$: "old": $\frac{\beta}{2}\vec{t}^T\vec{t}$

"new": $\frac{\beta}{2}(\vec{t}-\phi\vec{m}_N)^T(\vec{t}-\phi\vec{m}_N) + \frac{\alpha}{2}\vec{m}_N^T\vec{m}_N +$

$$+ \frac{1}{2}\vec{m}_N^T A \vec{m}_N = \frac{\beta}{2}\vec{t}^T\vec{t} + \frac{\beta}{2}(\phi\vec{m}_N)^T(\phi\vec{m}_N) -$$

$$-\underbrace{\beta(\phi\vec{m}_N)^T\vec{t}}_{\beta\vec{m}_N^T\phi^T\vec{t}} + \frac{\alpha}{2}\vec{m}_N^T\vec{m}_N + \frac{1}{2}\vec{m}_N^T A \vec{m}_N =$$

$$\underbrace{\beta\vec{m}_N^T\phi^T\vec{t}}_{\substack{\uparrow \\ \text{insert } AA^{-1}}} = \beta\vec{m}_N^T A \vec{m}_N$$

$$= \frac{\beta}{2}\vec{t}^T\vec{t} + \frac{\beta}{2}(\phi\vec{m}_N)^T(\phi\vec{m}_N) + \frac{\alpha}{2}\vec{m}_N^T\vec{m}_N -$$

$$- \frac{1}{2}\underbrace{\vec{m}_N^T A \vec{m}_N}_{\vec{m}_N^T \alpha \vec{m}_N + \beta(\phi\vec{m}_N)^T(\phi\vec{m}_N)} =$$

$$= \frac{\beta}{2}\vec{t}^T\vec{t} \quad \text{same!}$$

$$\underline{\underline{\quad}}$$

$\boxed{1'}$

Finally, $\theta(\vec{w}1)$:

"old": $-\beta \, (\varphi\vec{w})^T \vec{t}$

"new": $-\vec{m}_N^T A \vec{w} = -\beta \, \vec{t}^T \varphi (A^{-1})^T \overset{A^T}{"A"} \vec{w} =$

$= -\beta \, \vec{t}^T \varphi\vec{w} = -\beta \, (\varphi\vec{w})^T \vec{t} \quad "(A^T)^{-1}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ Same!

$A^T = \alpha I + \beta (\varphi^T \varphi)^T = \alpha I + \beta \, \varphi^T \varphi = A \, !$

So, all terms match between (3.79) & (3.80).

_____

Note also that

$$p(\vec{t} \mid \vec{w}, \beta) = \prod_{n=1}^{N} \mathcal{N}(t_n \mid \vec{w}^T \vec{\varphi}(\vec{x}_n), \beta^{-1}) =$$

$$= \mathcal{N}(\vec{t} \mid \underbrace{\varphi\vec{w}}_{\substack{N\text{-dim} \\ \text{vector}}}, \underbrace{\beta^{-1} I}_{\substack{N\times N \text{ diag.} \\ \text{cov. matrix}}})$$

Then

$$p(\vec{t}|\alpha,\beta) = \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \left(\frac{\alpha}{2\pi}\right)^{\frac{D+1}{2}} e^{-E(\vec{m}_N)} \times$$

$$\times \int d\vec{w}\, e^{-\frac{1}{2}(\vec{w}-\vec{m}_N)^T A (\vec{w}-\vec{m}_N)}$$

$$\underbrace{(2\pi)^{\frac{D+1}{2}}}_{\substack{\text{complexity} \\ \text{penalty}}} \underbrace{\frac{1}{|A|^{1/2}}}_{\det(A)}$$

recall that        13.82)

$$\swarrow\ E(\vec{m}_N) = \frac{\beta}{2}\|\vec{t}-\phi\vec{m}_N\|^2 + \\ + \frac{\alpha}{2}\vec{m}_N^T\vec{m}_N$$

$$\log p(\vec{t}|\alpha,\beta) = \frac{D+1}{2}\log\alpha + \frac{N}{2}\log\beta - \frac{N}{2}\log(2\pi) -$$

$$- E(\vec{m}_N) - \underbrace{\frac{1}{2}\log|A|}_{\text{complexity penalty}}$$

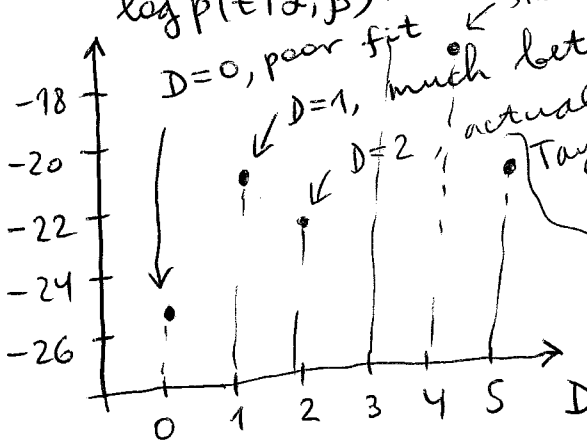$p(\vec{t}|\alpha,\beta)$ is basically the model evidence ⟨⇔⟩

⟨⇔⟩ $p(D|M_i)$

Recall the polynomial regression problem:

$$\text{fit}\quad y(x,\vec{w}) = \sum_{j=0}^{D} w_j\, x^j$$

to $\sin(2\pi x) + \text{noise}$.    $\underline{N=10}$

Choose $\alpha = 5\times10^{-3}$ and assume $\beta$ is known as well

$\log p(\vec{t}|\alpha,\beta)$ ← significant improvement    complexity penalty dominates

← small improvement,

D=0, poor fit

↙ D=1, much better fit

↙ D=2 ↑ actually worse since Taylor expansion of $\sin(2\pi x)$ has no even powers of $x$ →

an improvement in $E(\vec{m}_N)$ is offset by complexity penalty



Clear preference for  M=3

-6-

# Maximizing the evidence function

Consider maximizing $p(\vec{t}|\alpha,\beta)$ wrt $\alpha$.

Define $(\beta\,\phi^T\phi)\,\vec{u}_i = \lambda_i\,\vec{u}_i$ eigenvector eq'n

Note that $A$ has eigenvalues of $\alpha + \lambda_i$

Now, consider

$$\frac{\partial}{\partial\alpha}\log|A| = \frac{\partial}{\partial\alpha}\log\left\{\prod_{i=0}^{D}(\lambda_i+\alpha)\right\} =$$

$$= \sum_i \frac{1}{\lambda_i+\alpha}.$$

Then $\dfrac{\partial}{\partial\alpha}\log p(\vec{t}|\alpha,\beta) = \dfrac{D+1}{2\alpha} - \dfrac{1}{2}\vec{m}_N^T\vec{m}_N -$

$$-\frac{1}{2}\sum_i\frac{1}{\lambda_i+\alpha} = 0 \quad, \text{ or}$$

$$\alpha\,\vec{m}_N^T\vec{m}_N = \underbrace{(D+1)}_{\sum_i\frac{\lambda_i+\alpha}{\lambda_i+\alpha}} - \alpha\sum_i\frac{1}{\lambda_i+\alpha} \equiv \gamma,$$

where

$$\gamma = \sum_i \frac{\lambda_i}{\lambda_i+\alpha}.$$

So, $\hat{\alpha} = \dfrac{\gamma}{\vec{m}_N^T\vec{m}_N}$

implicit eq'n since both $\gamma$ & $\vec{m}_N$ depend on $\alpha$

Solve iteratively:
choose $\alpha \to$ compute $\gamma, \vec{m}_N \to$ update $\alpha$, etc.
Note that $\lambda_i$'s can be computed <u>once</u>.

Note that we assume

$$\frac{\partial}{\partial \alpha} \vec{m}_N \simeq 0, \quad \frac{\partial}{\partial \beta} \vec{m}_N \simeq 0.$$

Indeed, $\begin{cases} \vec{m}_N = \beta A^{-1} \varphi^T \vec{t}, \\ A = \alpha I + \beta \varphi^T \varphi. \end{cases}$

$\Downarrow$

$$\frac{\partial \vec{m}_N}{\partial \alpha} = \beta \underbrace{\frac{\partial A^{-1}}{\partial \alpha}} \varphi^T \vec{t} = \beta (A^{-1})^2 \varphi^T \vec{t}. \quad \leftarrow \text{suppressed by } A^{-1} \text{ compared to } \vec{m}_N$$

$$-A^{-1} \frac{\partial A}{\partial \alpha} A^{-1} = -(A^{-1})^2$$

$\overrightarrow{\text{derive using}}$

$$\frac{\partial}{\partial \alpha} \left[ A^{-1} A = I \right]$$

Further,

$$\frac{\partial \vec{m}_N}{\partial \beta} = \overbrace{+\beta \frac{\partial A^{-1}}{\partial \beta}}^{A^{-1}\varphi^T\vec{t} +} \varphi^T \vec{t} = \overbrace{-\beta A^{-1} \underbrace{\frac{\partial A}{\partial \beta}}_{\varphi^T\varphi} A^{-1}}^{A^{-1}\varphi^T\vec{t} -} \varphi^T \vec{t} =$$

$$= A^{-1} \varphi^T \vec{t} - A^{-1} \left[ A - \alpha I \right] A^{-1} \varphi^T \vec{t} =$$

$$= \alpha (A^{-1})^2 \varphi^T \vec{t}. \quad \leftarrow \text{suppressed by } A^{-1} \text{ compared to } \vec{m}_N$$

**3'**

Now, compute $\frac{\partial}{\partial \beta} \log p(\vec{t}|\alpha, \beta)$.

Note that $\lambda_i \sim \beta \implies \frac{\partial \lambda_i}{\partial \beta} = \frac{\lambda_i}{\beta}$.

Then $\frac{\partial}{\partial \beta} \log |A| = \frac{\partial}{\partial \beta} \sum_i \log(\lambda_i + \alpha) =$

$$= \frac{1}{\beta} \sum_i \frac{\lambda_i}{\lambda_i + \alpha} = \frac{\gamma}{\beta}, \quad \text{and}$$

$$\frac{\partial}{\partial \beta} \log p = \frac{N}{2\beta} - \frac{\gamma}{2\beta} - \frac{1}{2} \sum_{n=1}^{N} (t_n - \vec{m}_N^T \vec{\varphi}(\vec{x}_n))^2 = 0, \quad \text{or}$$

$$\frac{1}{\hat{\beta}} = \frac{1}{N-\gamma} \sum_{n=1}^{N} (t_n - \vec{m}_N^T \vec{\varphi}(\vec{x}_n))^2.$$

again, solve iteratively:

$$\beta \to \vec{m}_N, \gamma \to \text{update } \beta, \text{etc.}$$

So, now we have $\boxed{\hat{\alpha}, \hat{\beta}}$:

## predictive distribution

$$p(t|\vec{t}) = p(t|\vec{t}, \hat{\alpha}, \hat{\beta}) =$$

$$= \int d\vec{w}\, \underbrace{p(t|\vec{w}, \hat{\beta}) p(\vec{w}|\vec{t}, \hat{\alpha}, \hat{\beta})}_{\text{same as (3.57)}} = \mathcal{N}(t| \vec{m}_N^T \vec{\varphi}(\vec{x}), \quad \sigma_N^2(\vec{x})),$$

new value of $\vec{x}$

where $\begin{cases} \vec{m}_N = \hat{\beta} S_N \Phi^T \vec{t}, \quad S_N^{-1} = \hat{\alpha} I + \hat{\beta} \Phi^T \Phi, \\ \sigma_N^2(\vec{x}) = \beta^{-1} + \vec{\varphi}^T(\vec{x}) S_N \vec{\varphi}(\vec{x}). \end{cases}$

# Effective number of prms

Consider $\alpha = \dfrac{\gamma}{\overline{m}_N^T \overline{m}_N}$ again.

Recall that $\lambda_i$'s are eigenvalues of a positive definite matrix: $\lambda_i > 0$.

$$0 < \frac{\lambda_i}{\lambda_i + \alpha} \leq 1 \quad \Rightarrow \quad 0 < \gamma \leq D+1$$
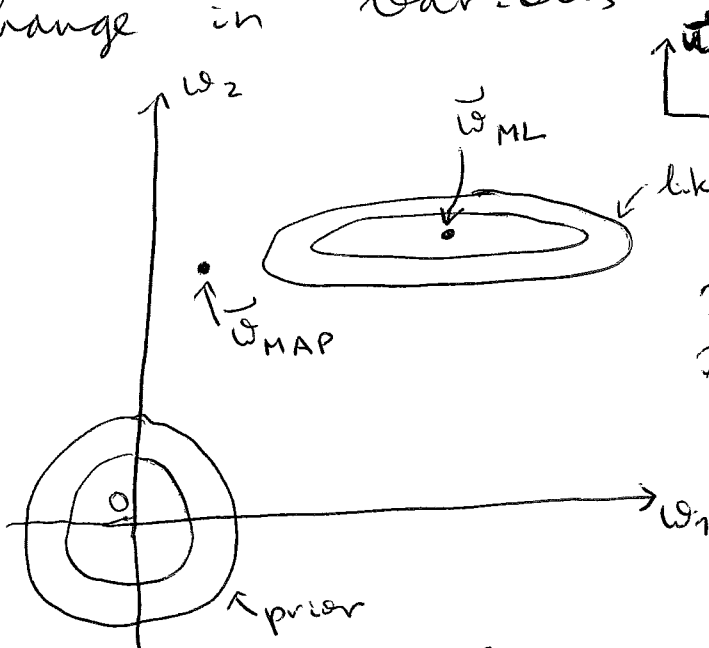
If $\lambda_i \gg \alpha \Rightarrow \dfrac{\lambda_i}{\lambda_i + \alpha} \approx 1$, and

$$\omega_i \approx \omega_i^{ML}.$$

Conversely, if $\lambda_i \ll \alpha \Rightarrow \dfrac{\lambda_i}{\lambda_i + \alpha} \approx 0,$

and $\omega_i \approx 0$.

Indeed, recall that $A = \alpha I + \beta \varphi^T \varphi$ is the curvature in the $\omega$ space since 

$$E(\vec{\omega}) = E(\overline{m}_N) + \frac{1}{2}(\vec{\omega} - \overline{m}_N)^T A (\vec{\omega} - \overline{m}_N).$$

Thus eigenvalues of $A$, $\lambda_i + \alpha$, will determine how quickly $\vec{\omega}$'s change in various directions.



Here,

$\lambda_1 \ll \alpha \Rightarrow \omega_{MAP,1} \approx 0$

$\lambda_2 \gg \alpha \Rightarrow \omega_{MAP,2} \approx \omega_{ML,2}$

Thus $\gamma$ determines the eff. # prms.

$\searrow$ $\gamma$ non-zero prms+
$+[(D+1)-\gamma]$ nearly zero prms

Note that

$$\frac{1}{\hat{\beta}} = \frac{1}{N-\gamma} \sum_{n=1}^{N} (t_n - \vec{m}_N^T \vec{\varphi}(\vec{x}_n))^2$$

$\underbrace{\phantom{\frac{1}{\hat{\beta}}}}$
"MAP" value

Recall that

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^{N} (t_n - \vec{w}_{ML}^T \vec{\varphi}(\vec{x}_n))^2$$

$N \to N-\gamma$ in the Bayesian result, correcting for the "ML bias".

---

If we consider the $N \gg D+1$ limit,

$\gamma = M$ since $\beta \varphi^T \varphi \uparrow$ as $N \uparrow$ due to a sum $\sum_{n=1}^{N}$. Thus $\lambda_i \uparrow$ as $N \uparrow$, $\forall i$

and $\frac{\lambda_i}{\lambda_i + d} \to 1 \implies \gamma \to M$.

In this case,

$$\left[ \quad \hat{d} \simeq \frac{M}{\vec{m}_N^T \vec{m}_N} \quad , \quad \hat{\beta}^{-1} \simeq \frac{1}{N} \sum_{n=1}^{N} (t_n - \vec{m}_N^T \vec{\varphi}(\vec{x}_n))^2 \right]$$

still iterative but easier since $\gamma$ does not need to be re-computed.

**Ex.:**

(β known
for simplicity)



$\gamma$ ($\approx D+1$ on the left,
$\to 0$ on the right)

$\leftarrow 2\tilde{m}_N^T \tilde{m}_N$

$\leftarrow \log p(\vec{t}|\alpha,\beta)$,
max at $\hat{\alpha}$

$-5 \qquad 0 \qquad \hat{\alpha} \qquad 5 \quad \log\alpha$