

Lecture 3

Define

$$h(\vec{x}) = \int dt t p(t|\vec{x}) = E[t|\vec{x}]$$

optimal prediction under squared loss f'n

Then $E[L] = \int d\vec{x} p(\vec{x}) [y(\vec{x}) - h(\vec{x})]^2 +$

$$+ \int dt d\vec{x} p(\vec{x}, t) [h(\vec{x}) - t]^2$$

indep. of $y(\vec{x})$, reflects intrinsic noise/scatter in the data

Imagine modeling $h(\vec{x})$ with $y(\vec{x}, \vec{w})$.
Consider multiple (K) datasets of size N .
For each dataset k , we can obtain

$$y(\vec{x}, \vec{w}_{ML,k}) \equiv y(\vec{x}; \mathcal{D}_k)$$

↑ particular dataset

Now, consider

$$(y(\vec{x}; \mathcal{D}_k) - h(\vec{x}))^2 = (y(\vec{x}; \mathcal{D}_k) - E(y(\vec{x}; \mathcal{D}_k)) + E(y(\vec{x}; \mathcal{D}_k)) - h(\vec{x}))^2 \Rightarrow$$

$$E(y(\vec{x}; \mathcal{D}_k)) = \frac{1}{K} \sum_{k=1}^K y(\vec{x}; \mathcal{D}_k)$$

$$\Rightarrow E[(y(\vec{x}; \mathcal{D}_k) - h(\vec{x}))^2] =$$

$$= E[(y(\vec{x}; \mathcal{D}_k) - E(y(\vec{x}; \mathcal{D}_k)))^2] +$$

$$+ E[(E(y(\vec{x}; \mathcal{D}_k)) - h(\vec{x}))^2] +$$

$$+ 2 E[(y(\vec{x}; \mathcal{D}_k) - E(y(\vec{x}; \mathcal{D}_k))) (E(y(\vec{x}; \mathcal{D}_k)) - h(\vec{x}))]$$

since
indep. of the E-operator

$$E(y) - E(E(y)) = E(y) - E(y) = 0$$

$$\boxed{E[(E(y) - h(\vec{x}))^2] + E[(y - E(y))^2]}$$

" (E(y) - h(\vec{x}))^2 "
variance

bias²

So, expected loss = (bias)² + variance + noise
 ↑
 to be minimized

Expect "rigid" models to have high bias / low variance & "flexible" models to have low bias / high variance

above, a single \vec{x} was considered. Over all values of \vec{x} ,

$$(bias)^2 = \int d\vec{x} p(\vec{x}) [E[y(\vec{x}; \mathcal{D}_k)] - h(\vec{x})]^2$$

$$variance = \int d\vec{x} p(\vec{x}) E[(y(\vec{x}; \mathcal{D}_k) - E(y(\vec{x}; \mathcal{D}_k)))^2]$$

$$noise = \int dt d\vec{x} p(\vec{x}, t) [h(\vec{x}) - t]^2$$

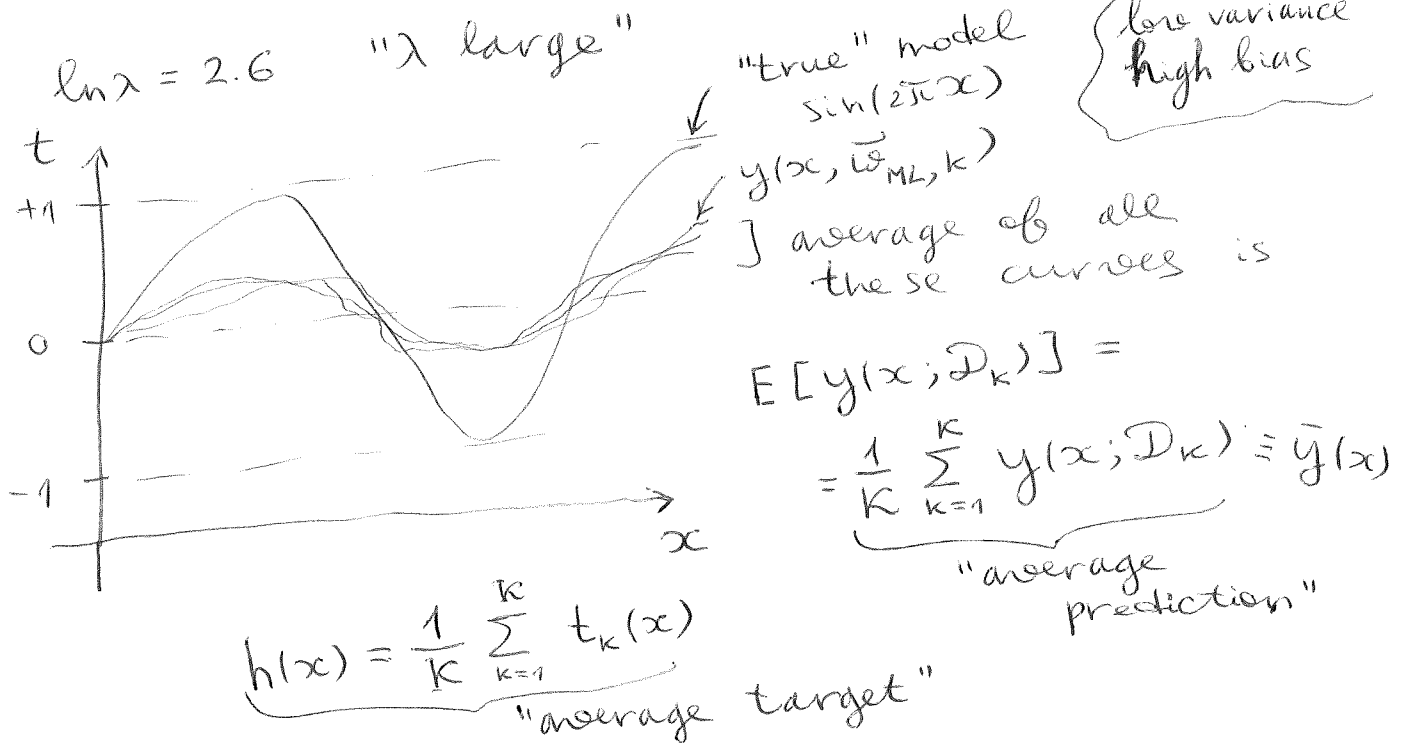
Example

Consider $K=100$ data sets, each with $N=25$ points, sampled from ~~the~~ $\sin(2\pi x) + \text{noise}$.

To each data set, we fit a model with 24 Gaussian basis functions by minimizing

$$\tilde{E}(\vec{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \vec{w}^T \vec{\phi}(\vec{x}_n))^2 + \frac{\lambda}{2} \vec{w}^T \vec{w}$$

Consider several values of λ :



Finally,

$$(\text{bias})^2 = \frac{1}{N} \sum_{n=1}^N [\bar{y}(x_n) - h(x_n)]^2$$

$$\text{variance} = \frac{1}{N} \sum_{n=1}^N \frac{1}{K} \sum_{k=1}^K [y(x_n, \vec{w}_{ML,k}) - \bar{y}(x_n)]^2$$

$$\text{noise} = \frac{1}{N} \sum_{n=1}^N \frac{1}{K} \sum_{k=1}^K (h(x_n) - t_k(x_n))^2$$

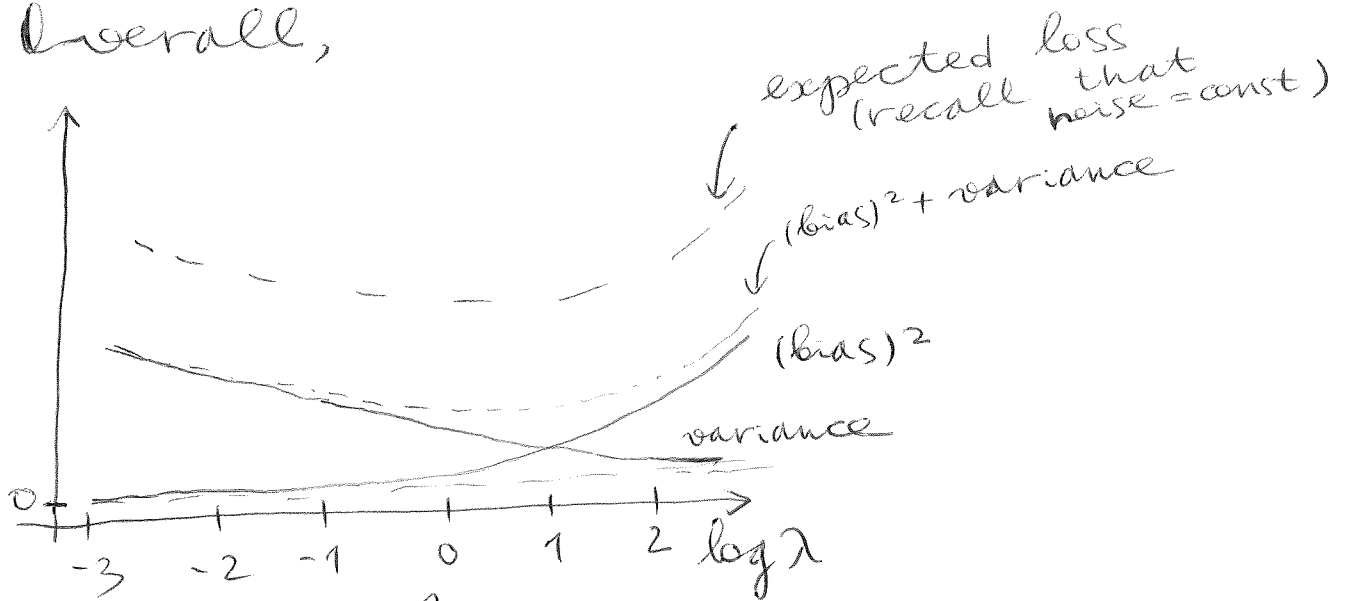
Now, consider $\ln \lambda = -0.31$ "medium"



If we try $\ln \lambda = -2.4$ "low"

{ very low bias
high variance }

Overall,



choose λ to minimize expected loss and either refit the final model on all data or just average K fit with the best value of λ
{ separate models }

Bayesian linear regression

Assume $p(\vec{w} | \alpha) = \mathcal{N}(\vec{w} | \vec{0}, \alpha^{-1} I)$
prior

Then the posterior is gaussian as well, since the likelihood

$$p(\vec{t} | \vec{w}) \propto \prod_{n=1}^N \mathcal{L}^{-\frac{\beta}{2} \sum_{n=1}^N (t_n - \vec{w}^T \vec{\phi}(\vec{x}_n))^2}$$

[assume that β is known for now] is gaussian & we can complete the square:

$$p(\vec{w} | \vec{t}) \sim \underbrace{\mathcal{L}^{-\frac{\beta}{2} \sum_{n=1}^N (t_n - \vec{w}^T \vec{\phi}(\vec{x}_n))^2}}_{\text{likelihood}} \underbrace{\mathcal{L}^{-\frac{\alpha}{2} \vec{w}^T \vec{w}}}_{\text{prior}}$$

Further,

$$p(\vec{w} | \vec{t}) = \mathcal{N}(\vec{w} | \vec{m}_N, \Sigma_N)$$

\uparrow \uparrow
 $D+1$ $(D+1) \times (D+1)$
 vector covariance matrix

$N = \# \text{ datapoints}$

Recall that

$$(*) \quad \mathcal{N}(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \mathcal{L}^{-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})}$$

\downarrow \downarrow
 mean vector covariance matrix
 D -dim gaussian $\det(\Sigma)$

\sim infinitely broad prior

$\alpha \rightarrow 0$:

$$\vec{w}_{\text{MAP}} = \vec{m}_N \quad \left\{ \begin{array}{l} \vec{m}_N = \beta S_N \Phi^T \vec{t} \\ S_N^{-1} = \alpha I + \beta \Phi^T \Phi \end{array} \right. \quad \vec{m}_N \rightarrow \underbrace{(\Phi^T \Phi)^{-1} \Phi^T \vec{t}}_{\vec{w}_{\text{ML}}}$$

(mean = mode for gaussians)

(2.113) - (2.117):

$$\begin{cases} p(\vec{x}) = \mathcal{N}(\vec{x} | \vec{\mu}', \Lambda^{-1}) \\ p(\vec{y} | \vec{x}) = \mathcal{N}(\vec{y} | A\vec{x} + \vec{b}, L^{-1}) \end{cases} \Rightarrow p(\vec{x} | \vec{y}) = \underbrace{\mathcal{N}(\vec{x} | \vec{m})}_{= \mathcal{N}(\vec{x} | S[A^T L (\vec{y} - \vec{b}) + \Lambda \vec{\mu}'], S)}, \text{ where}$$

$$S = (\Lambda + A^T L A)^{-1}$$

Here, $\begin{cases} p(\vec{w} | \alpha) = \mathcal{N}(\vec{w} | \vec{0}, \alpha^{-1} I) & \text{prior} \\ p(\vec{t} | \vec{w}) = \mathcal{N}(\vec{t} | \Phi \vec{w}, \beta^{-1} I) & \text{likelihood} \end{cases}$

↑
Indeed,

$$p(\vec{t} | \vec{w}) \sim \ell^{-\frac{1}{2}} \left[(t_1 - \underbrace{\sum_j w_j \phi_j(\vec{x}_1)}_{\text{sum over } j}) \right] \beta (t_1 - \sum_j w_j \phi_j(\vec{x}_1)) + \dots + (t_N - \sum_j w_j \phi_j(\vec{x}_N)) \beta (t_N - \sum_j w_j \phi_j(\vec{x}_N))$$

clearly, in eq'n (*) $\Sigma^{-1} = \beta I$ and

$$\underbrace{\vec{\mu}}_N = \begin{pmatrix} w_j \phi_j(\vec{x}_1) \\ \vdots \\ w_j \phi_j(\vec{x}_N) \end{pmatrix} = \begin{pmatrix} \phi_{1j} w_j \\ \vdots \\ \phi_{Nj} w_j \end{pmatrix} = \Phi \vec{w}$$

$$N \left\{ \underbrace{\begin{pmatrix} \Phi \\ \vec{0} \end{pmatrix}}_{D+1} \right\}_{D+1} = \underbrace{\vec{\mu}}_N, \text{ as expected}$$

So, $\vec{x} \rightarrow \vec{w} \Rightarrow \begin{cases} \vec{\mu}' = 0, \\ \Lambda^{-1} = \alpha^{-1} I \end{cases}$

$$\vec{y} \rightarrow \vec{t} \Rightarrow \begin{cases} A = \Phi, \vec{b} = 0 \\ L^{-1} = \beta^{-1} I \end{cases}$$

Finally, $S^{-1} = \alpha I + \beta \Phi^T \Phi$ and $\vec{m} = S \Phi^T \beta \vec{t} = \beta S \Phi^T \vec{t}$.
So, $S = S_N$ & $\vec{m} = \vec{m}_N$ from opposite page

We argued before that maximizing $p(\vec{w} | \vec{E})$ is like minimizing $\tilde{E}(\vec{w})$ with $\lambda = \alpha/\beta$.

Example ^{consider} $x \rightarrow t$ (both 1D)
_{input target}

$$y(x, \vec{w}) = w_0 + w_1 x \quad \text{linear model}$$

Synthetic data:

$$t = f(x, \vec{a}) + \text{noise}, \text{ where}$$

$$f(x, \vec{a}) = a_0 + a_1 x \quad \begin{cases} a_0 = -0.3 \\ a_1 = 0.5 \end{cases}$$

$$\text{noise} = \mathcal{N}(\mu=0, \sigma^2), \quad \sigma = 0.2$$

Sampling: choose x_n uniformly in the $[1, 1]$ range \Rightarrow evaluate

$$f(x_n, \vec{a}) \Rightarrow \text{generate value } \hat{t}_n \Rightarrow \text{get } t_n = f(x_n, \vec{a}) + \hat{t}_n$$

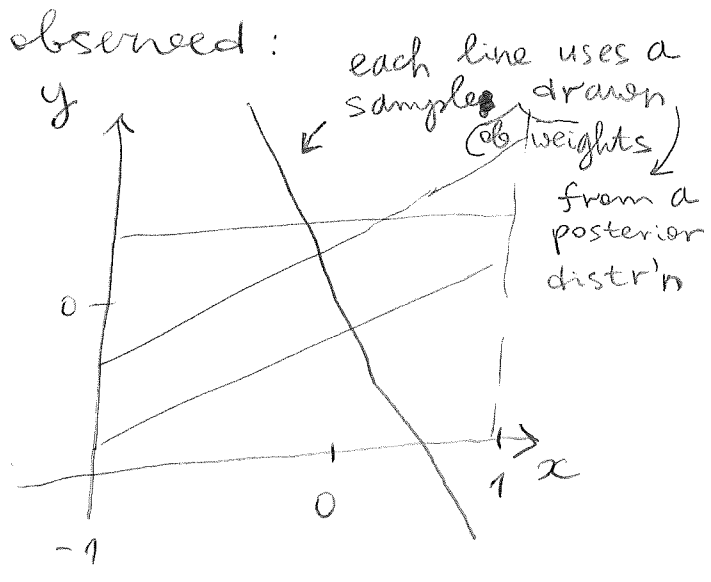
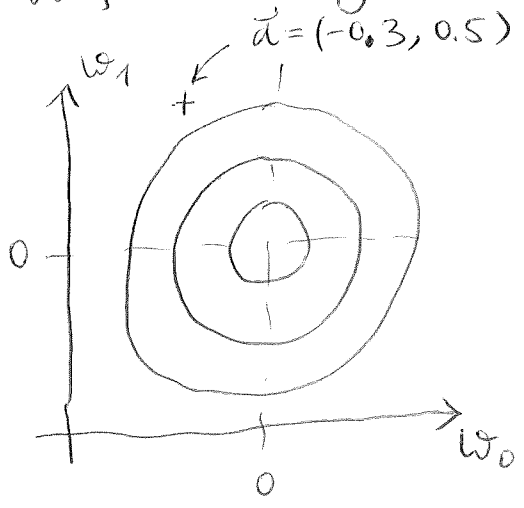
$$n = 1, \dots, N$$

Assume $\beta = \frac{1}{\sigma^2} = 25$ is known exactly.

Choose $\lambda = 2.0$

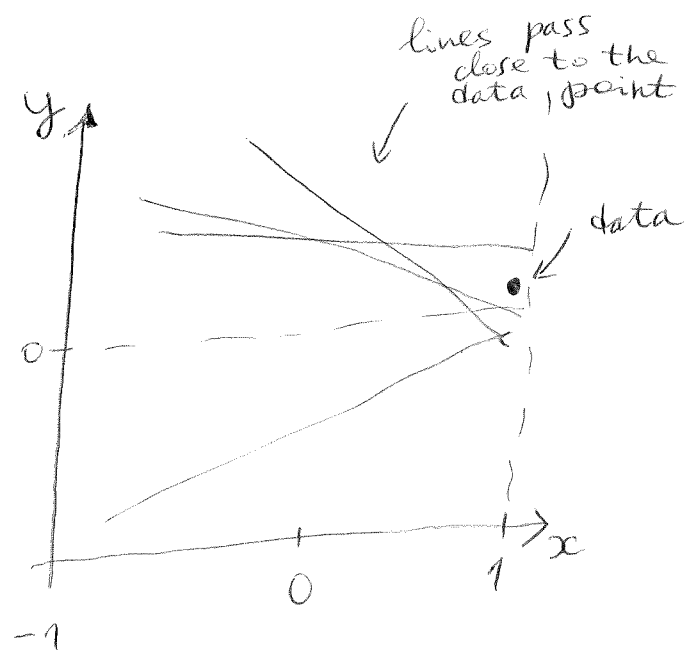
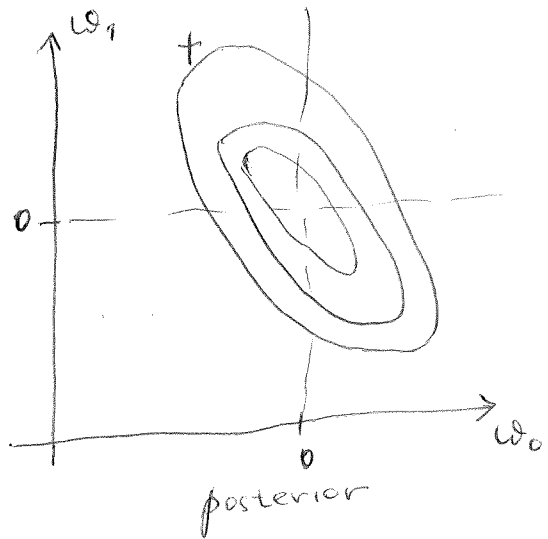
these will be relaxed later

Before any data is observed:

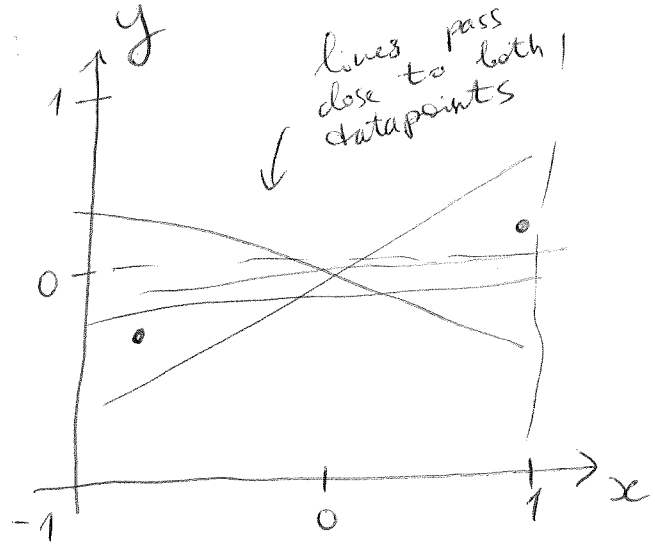
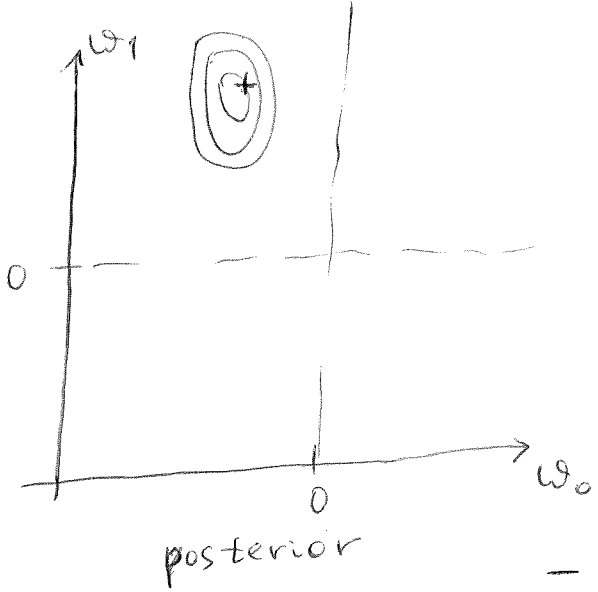


posterior = prior

One data point:



two data points:



Many data points:

