

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
$$E[L] = \int d\vec{x} p(\vec{x}) \underbrace{[y(\vec{x}) - h(\vec{x})]^2}_{\substack{\uparrow \\ \text{model}}} + \int dt d\vec{x} p(\vec{x}, t) [h(\vec{x}) - t]^2$$

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

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

\uparrow 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 $\underbrace{E(y) - E(E(y)) = E(y) - E(y) = 0}$ indep. of the E-operator

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

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

$$\underbrace{(E(y) - h(\vec{x}))^2}_{\text{bias}^2} \quad \underbrace{(y - E(y))^2}_{\text{variance}}$$

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

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

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

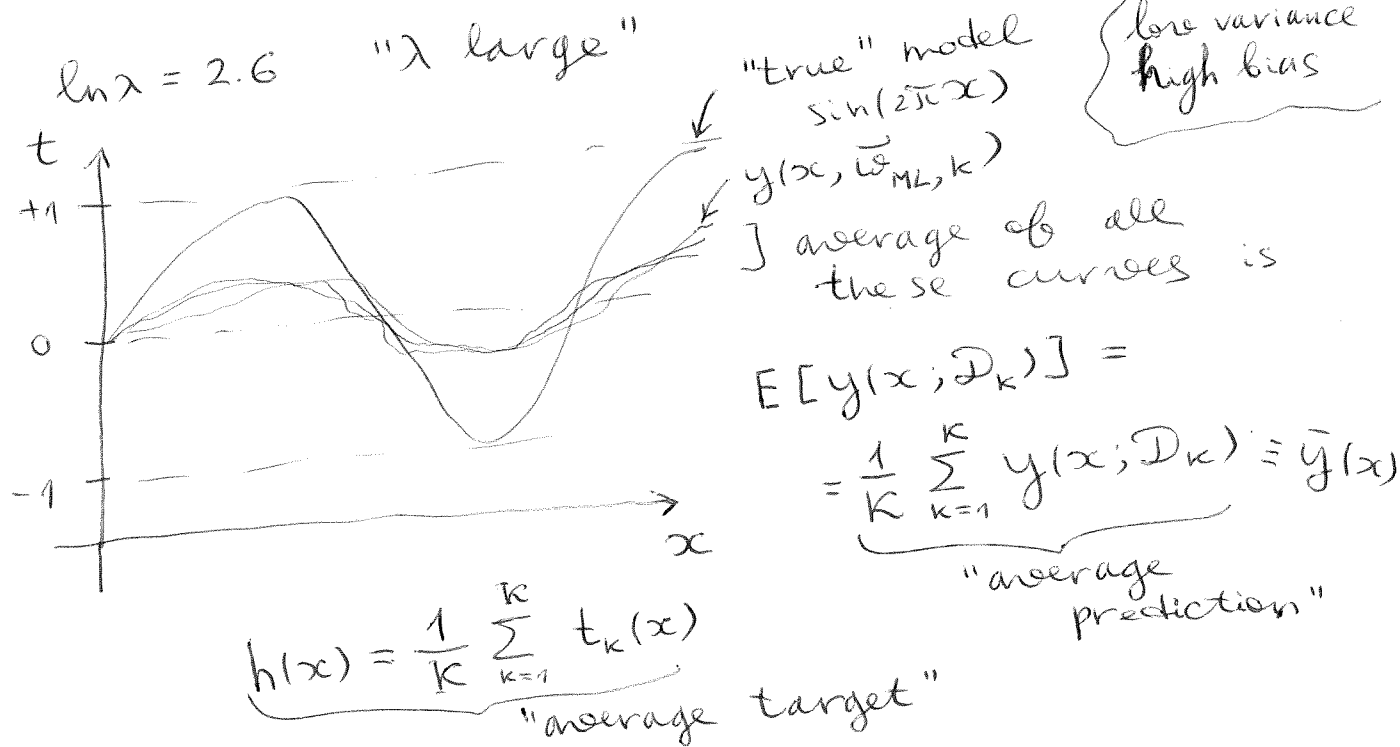
Example

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

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,

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

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

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

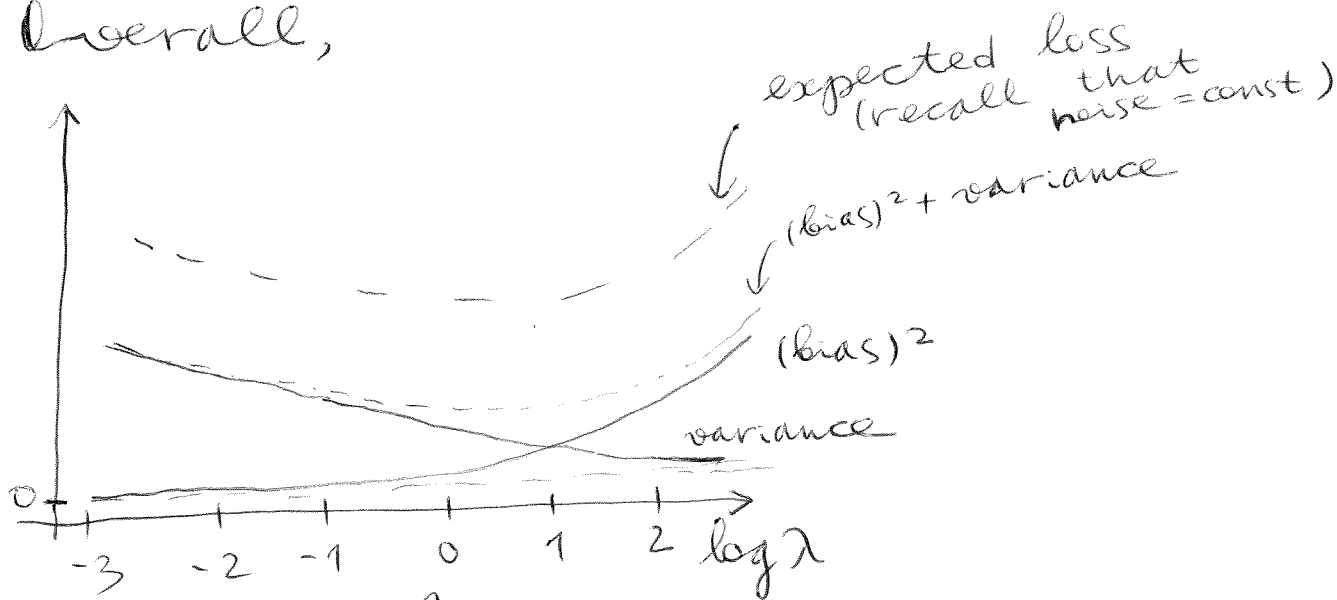
Now, consider $\ln \lambda = -0.31$ "medium"



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

Bayesian linear regression

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

Then the posterior is gaussian as well, since the likelihood

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

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

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

Further,

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

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

$N = \#$ datapoints

Recall that

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

D-dim gaussian det(Σ)

In our case,

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

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

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

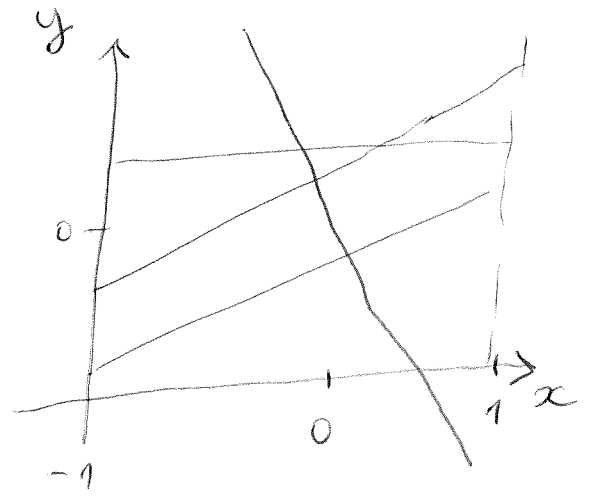
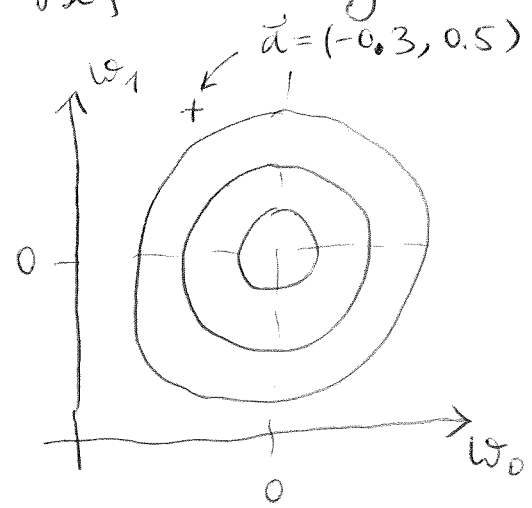
$$n = 1, \dots, N$$

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

Choose $\lambda = 2.0$

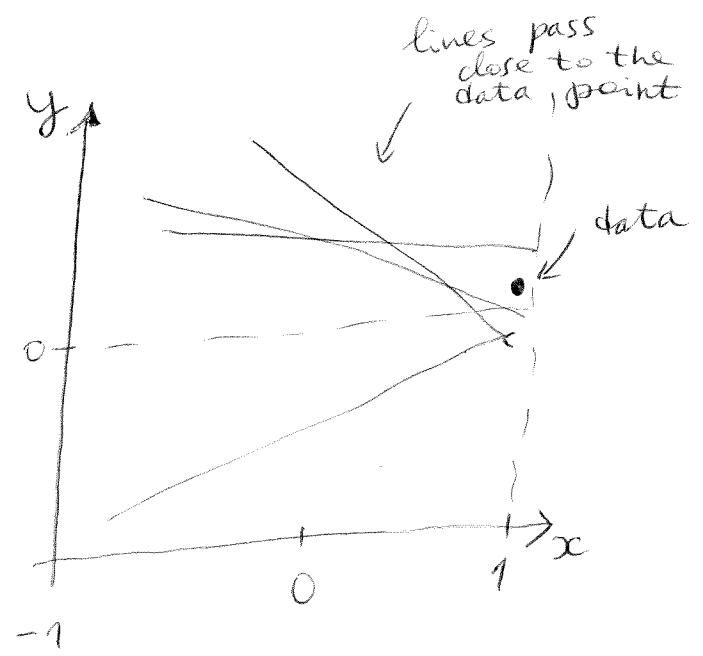
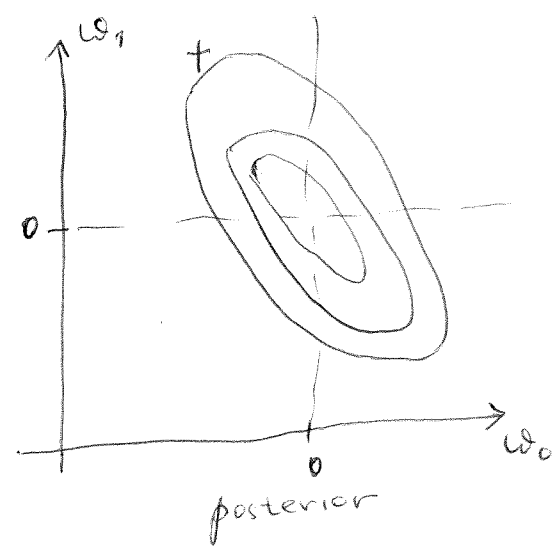
these will be relaxed later

Before any data is observed:

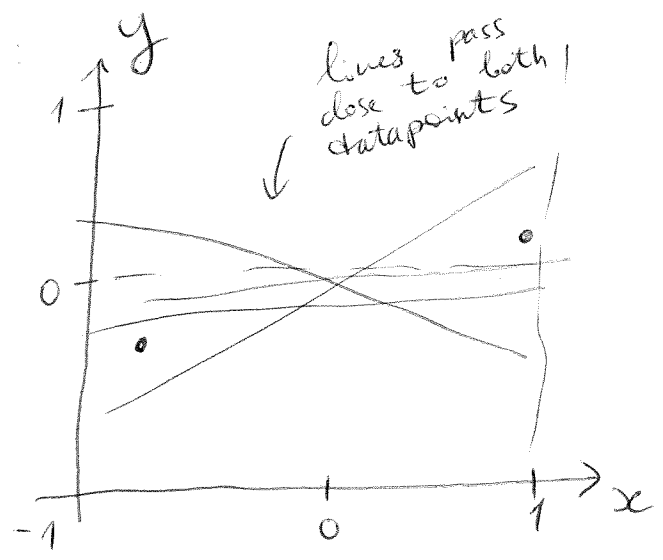
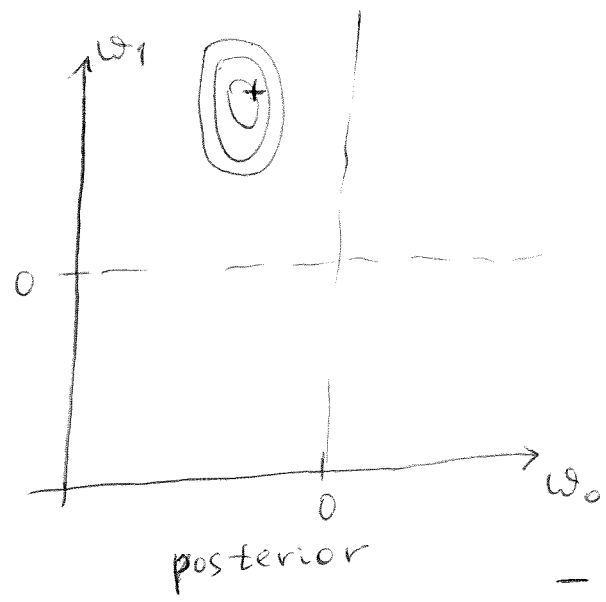


posterior = prior

One data point:



two data points:



Many data points:

