



Institute for Advanced Studies
in Basic Sciences
Gava Zang, Zanjan, Iran

Regression with Linear Models

By Dr. Parvin Razzaghi (p.razzaghi@iasbs.ac.ir)

Winter 2024





Slides credit

- Professor:
 - Dr. Parvin Razzaghi
- Original Slides:
 - Amir-massoud Farahman
 - Eric Eaton
- Link:
 - <https://B2n.ir/n25709>
- University:
 - University of Toronto and Vector Institute
- Acknowledgment:
 - Credit for slides goes to many members of the ML Group at the U of T, and beyond, including (recent past): Roger Grosse, Murat Erdogdu, Richard Zemel, Juan Felipe Carrasquilla, Emad Andrews.
 - For any comments or suggestions, please contact: p.razzaghi@iasbs.ac.ir



ℓ_2 or (L^2) Regularization

- The regularized cost function

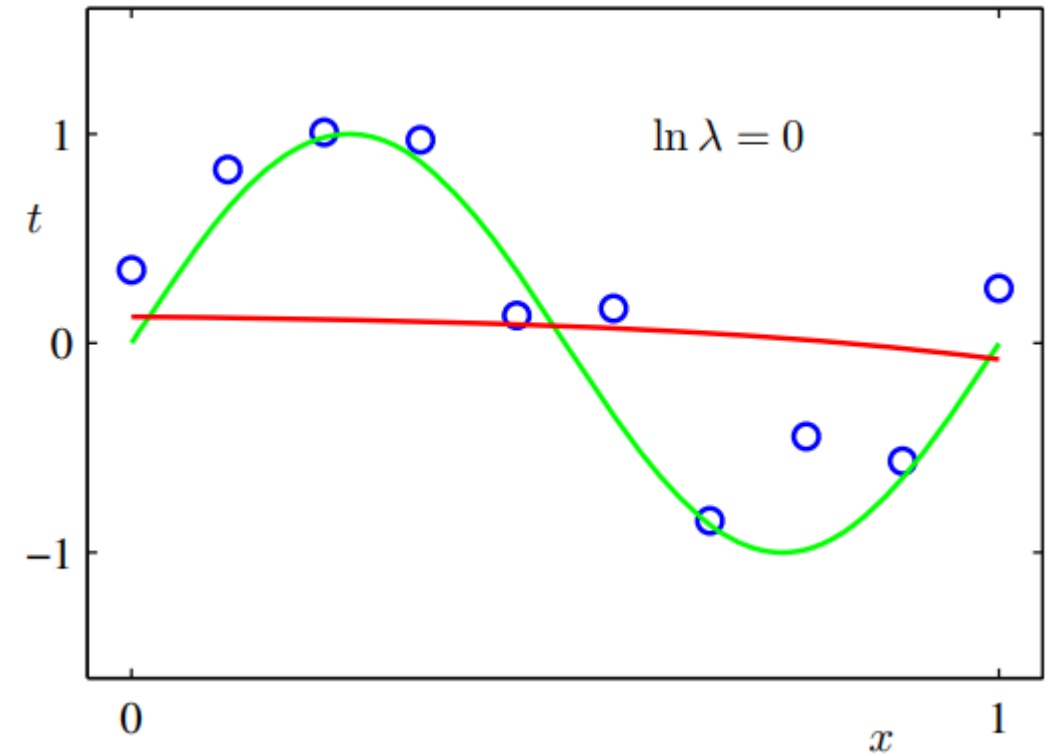
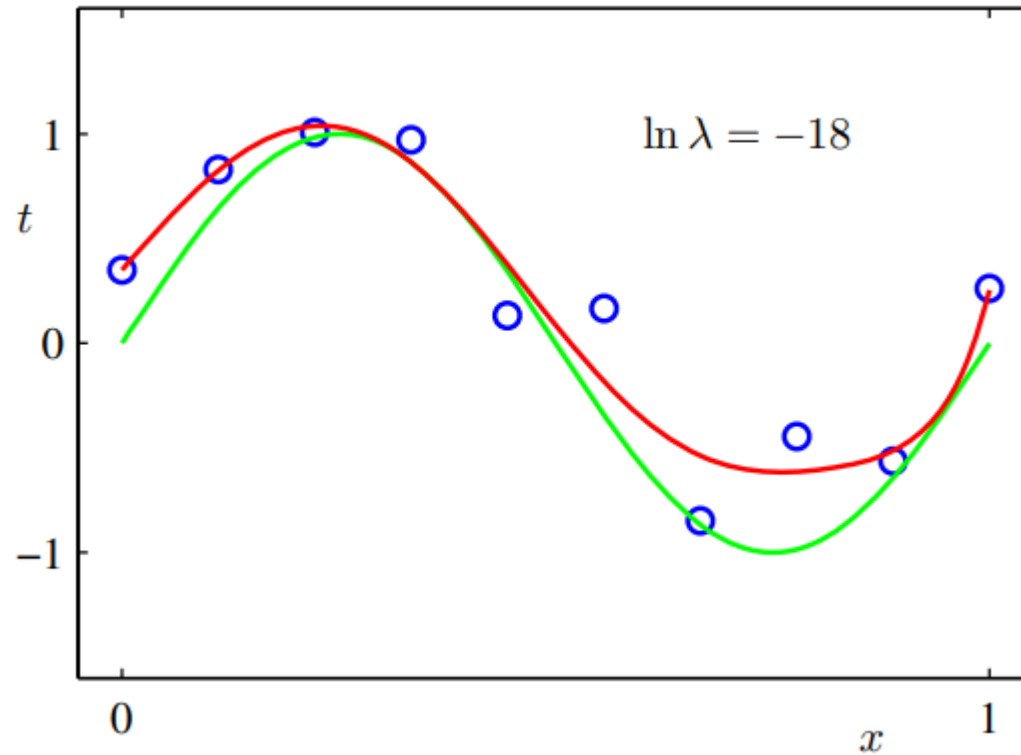
$$\mathcal{J}_{reg}(w) = \mathcal{J}(w) + \lambda \mathcal{R}(w) = \mathcal{J}(w) + \frac{\lambda}{2} \sum_j w_j^2.$$

- The basic idea is that “simpler” functions have weights w with smaller ℓ_2 -norm and we prefer them to functions with larger ℓ_2 -norm.
 - **Intuition:** Large weights makes the function f have more abrupt changes as a function of the input x ; it will be less smooth.
- If you fit training data poorly, \mathcal{J} is large. If the fitted weights have high values, \mathcal{R} is large.
- Large λ penalizes weight values more.
- Here λ is a hyperparameter that we can tune with a validation set.



Regularized linear regression

- $M=9$

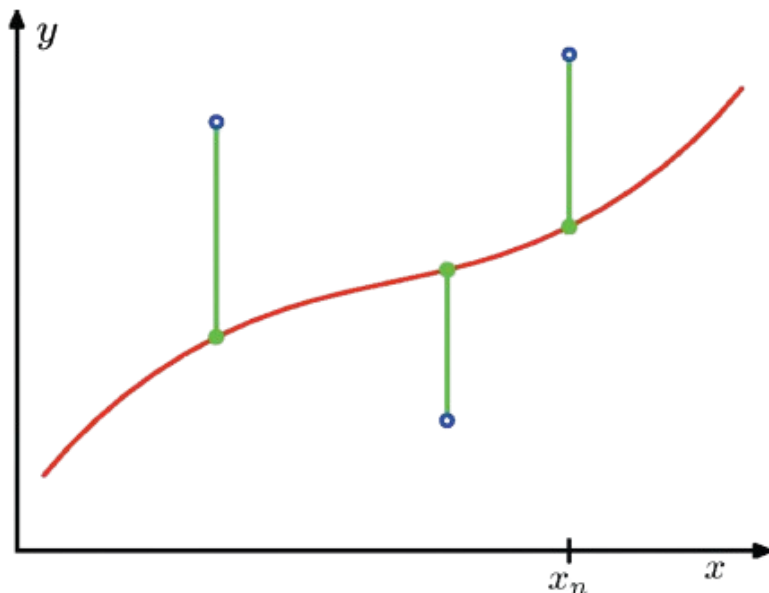




Probabilistic Interpretation of the Squared Error

- For the least squares: we minimize the sum of the squares of the errors between the predictions for each data point $\mathbf{x}^{(i)}$ and the corresponding target value $t^{(i)}$, i.e.,

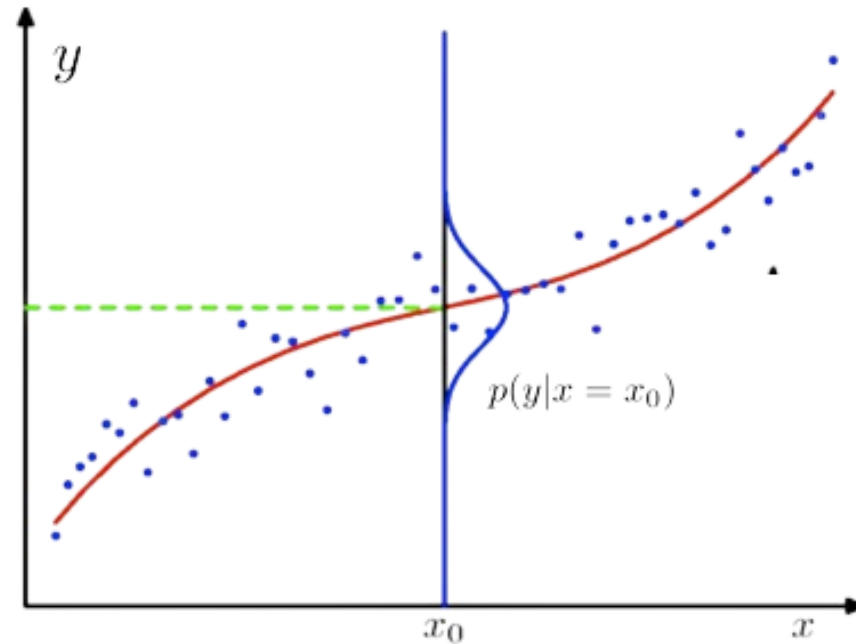
$$\underset{(\mathbf{w}, \mathbf{w}_0)}{\text{minimize}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} + b - t^{(i)})^2$$



- $t \approx \mathbf{x}^T \mathbf{w} + b, (\mathbf{w}, b) \in \mathbb{R}^D \times \mathbb{R}$
- We measure the quality of the fit using the squared error loss. Why?
- Even though the squared error loss is intuitive, we did not justify it.
- We provide a probabilistic perspective here.
- There are other justifications too; we get to them in the Bias-Variance decomposition lecture.



Probabilistic Interpretation of the Squared Error



- Suppose that our model arose from a statistical model ($b = 0$ for simplicity):

$$y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \varepsilon^{(i)}$$

where $\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$ is independent of the input $\mathbf{x}^{(i)}$.

- Thus, $y^{(i)} | \mathbf{x}^{(i)} \sim p(y | \mathbf{x}^{(i)}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^T \mathbf{x}^{(i)}, \sigma^2)$.



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

- Suppose that the input data $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$ are given and the outputs are independently drawn from

$$t^{(i)} \sim p(y|\mathbf{x}^{(i)}, \mathbf{w}).$$

- with an unknown parameter \mathbf{w} . So the dataset is

$$D = \{(\mathbf{x}^{(1)}, t^{(1)}), \dots, (\mathbf{x}^{(N)}, t^{(N)})\}.$$



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

- The **likelihood** function is $\Pr(D|\mathbf{w})$.
- The **maximum likelihood estimation (MLE)** is based on the “**principle**” suggesting that we have to find a parameter $\hat{\mathbf{w}}$ that maximizes the likelihood, i.e.,

$$\hat{\mathbf{w}} \leftarrow \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(D|\mathbf{w}) .$$

Maximum likelihood estimation: after observing the data samples $(\mathbf{x}^{(t)}, t^{(i)})$ for $i = 1, 2, \dots, N$, we should choose \mathbf{w} that maximizes the likelihood.



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

- For independent samples, the likelihood function of samples D is the product of their likelihoods

$$p(t^{(1)}, t^{(2)}, \dots, t^{(N)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}, w) = \prod_{i=1}^N p(t^{(i)}, \mathbf{x}^{(i)}, w) = L(w).$$

- Product of N terms is not easy to minimize.
- Taking log reduces it to a sum. Two objectives are equivalent since log is strictly increasing.
- Maximizing the likelihood is equivalent to minimizing the **negative log-likelihood**:



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

- Maximizing the likelihood is equivalent to minimizing the **negative log-likelihood**:

$$\ell(w) = -\log L(w) = -\log \prod_{i=1}^N p(t^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) = -\sum_{i=1}^n \log p(t^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}).$$



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

Maximum Likelihood Estimator (MLE)

After observing $\mathbf{z}^{(i)} = (\mathbf{x}^{(i)}, t^{(i)})$ for $i = 1, \dots, N$ independent and identically distributed (i.i.d.) samples from $p(\mathbf{z}, \mathbf{w})$, MLE is

$$\mathbf{w}^{MLE} = \underset{\mathbf{w}}{\operatorname{argmin}} l(\mathbf{w}) = - \sum_{i=1}^N \log p(t^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}).$$



Probabilistic Interpretation of the Squared Error: From MLE to Squared Error

- Suppose that our model arose from a statistical model:

$$y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \varepsilon^{(i)}$$

where $\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$ is independent of anything else.

- $p(t^{(i)}, \mathbf{x}^{(i)}, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2\right\}$
- $\log p(t^{(i)}, \mathbf{x}^{(i)}, w) = \frac{1}{2\sigma^2} (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 - \log(\sqrt{2\pi\sigma^2})$



Probabilistic Interpretation of the Squared Error: From MLE to Squared Error

- The MLE solution is

$$\mathbf{w}^{MLE} = \underset{\mathbf{w}}{\operatorname{argmin}} \mathcal{L}(\mathbf{w}) = \frac{1}{2\sigma^2} \sum_{i=1}^N (t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 + \mathcal{C}$$

- As \mathcal{C} and σ do not depend on \mathbf{w} , they do not contribute to the minimization.

$\mathbf{w}^{MLE} = \mathbf{w}^{LS}$ when we work with Gaussian densities.



Probabilistic Interpretation of the Squared Error: From MLE to Squared Error

- Suppose that our model arose from a statistical model:

$$y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \varepsilon^{(i)}$$

where $\varepsilon^{(i)}$ comes from the **Laplace** distribution, that is, the distribution of $\varepsilon^{(i)}$ has density.

$$\frac{1}{2b} \exp \left(-\frac{|y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)}|}{2b} \right)$$



Probabilistic Interpretation of the Squared Error: From MLE to Squared Error

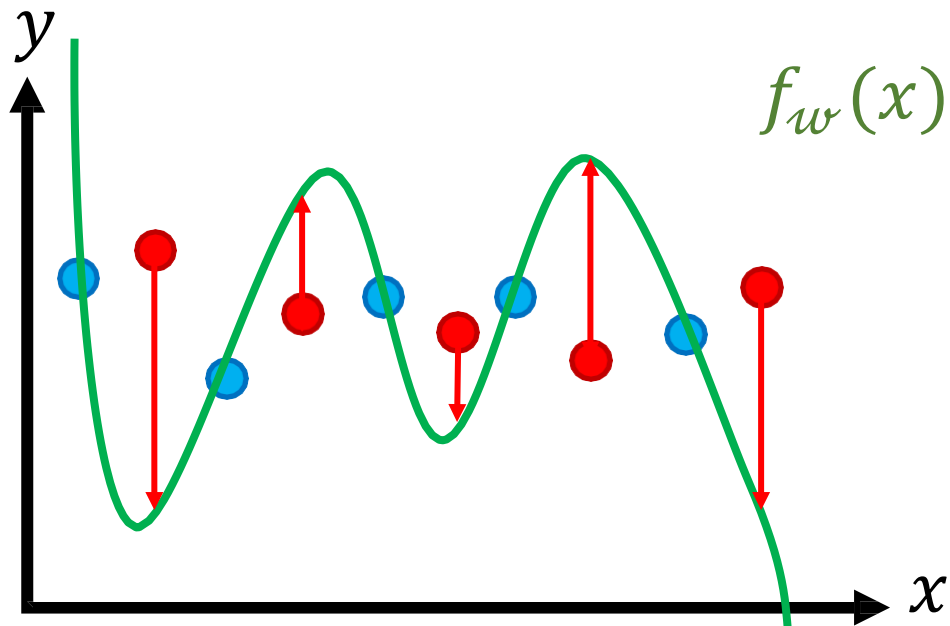
- **Q:** What is the loss in the MLE?
 - Choice 1: $\frac{1}{N} \sum_{i=1}^N |t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)}|^{1/2}$
 - Choice 2: $\frac{1}{N} \sum_{i=1}^N (t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})$
 - Choice 3: $\frac{1}{N} \sum_{i=1}^N |t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)}|$
 - Choice 4: $\frac{1}{N} \left| \sum_{i=1}^N t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)} \right|$
- **Q:** Can you think of an application area with non-Gaussian probabilistic model?



Bias-Variance Tradeoff

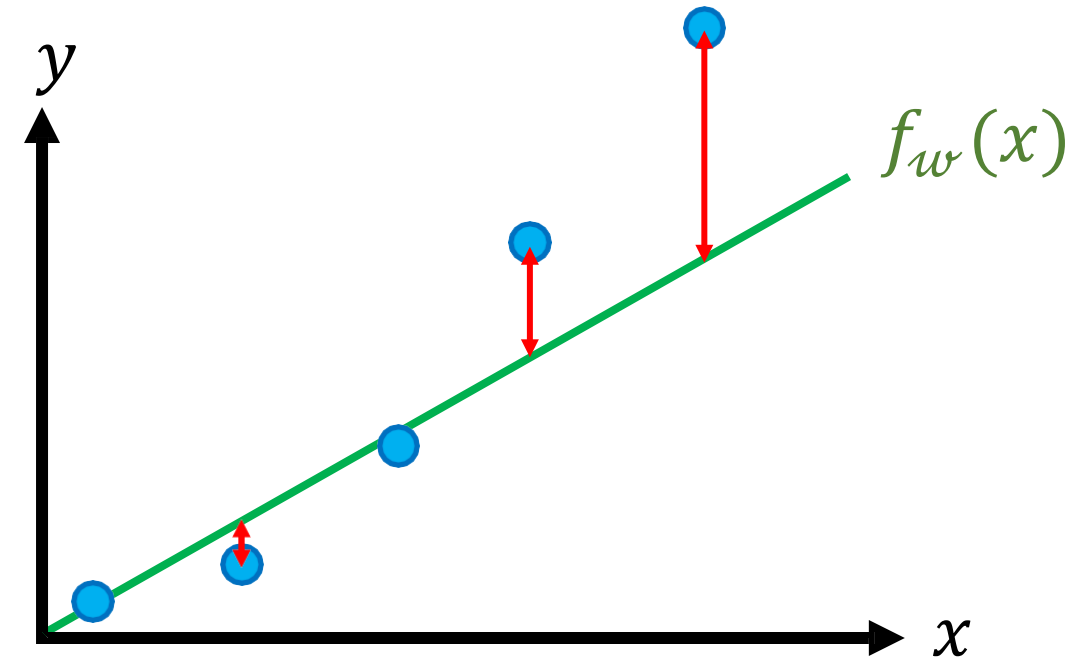
- **Overfitting (high variance)**

- High capacity model capable of fitting complex data
- Insufficient data to constrain it



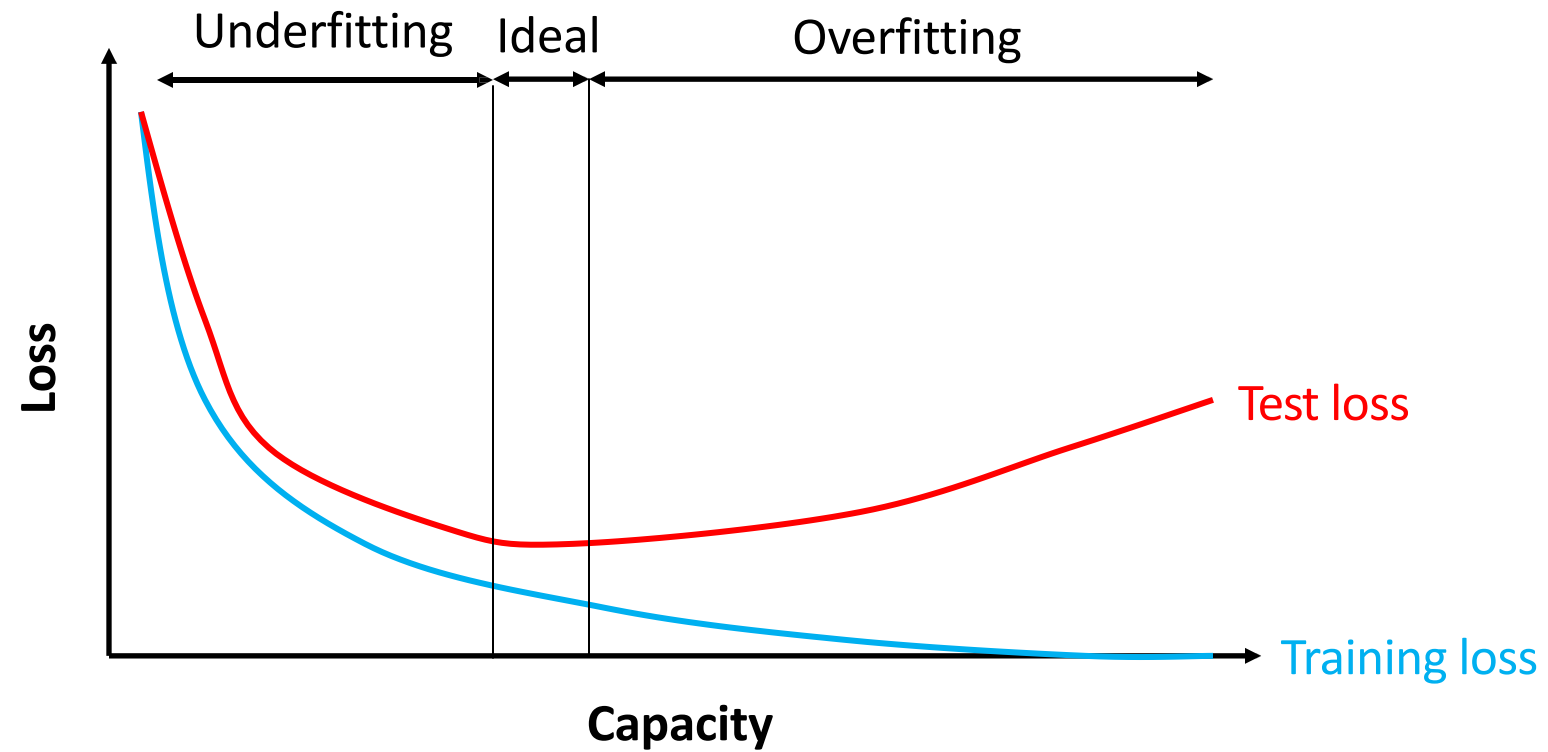
- **Underfitting (high bias)**

- Low capacity model that can only fit simple data
- Sufficient data but poor fit





Bias-Variance Tradeoff





The Bias-Variance Decomposition₍₁₎

- Recall the *expected squared loss*,

$$\mathbb{E}[L] = \int \{y(x) - h(x)\}^2 p(x) dx + \underbrace{\iint \{h(x) - t\}^2 P(x, t) dx dt}_{\text{noise}}$$

where

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

The second term of $E[L]$ corresponds to the noise inherent in the random variable t .

Q: What about the first term?



The Bias-Variance Decomposition₍₂₎

- Suppose we were given multiple data sets, each of size N . Any particular data set, \mathcal{D} , will give a particular function $y(x; \mathcal{D})$. We then have

$$\begin{aligned} & \{y(x; \mathcal{D}) - h(x)\}^2 \\ &= \{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\}^2 \\ &= \{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\}^2 \\ &+ 2\{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})]\} \{\mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\} \end{aligned}$$



The Bias-Variance Decomposition₍₃₎

- Taking the expectation over D yields

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}}[\{y(x; \mathcal{D}) - h(x)\}^2] \\ &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}}[\{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})]\}^2]}_{\text{variance}} \end{aligned}$$



The Bias-Variance Decomposition₍₄₎

- Thus we can write

$$\textit{expected loss} = (\textit{bias})^2 + \textit{variance} + \textit{noise}$$

where

$$(\textit{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\}^2 p(x) dx$$

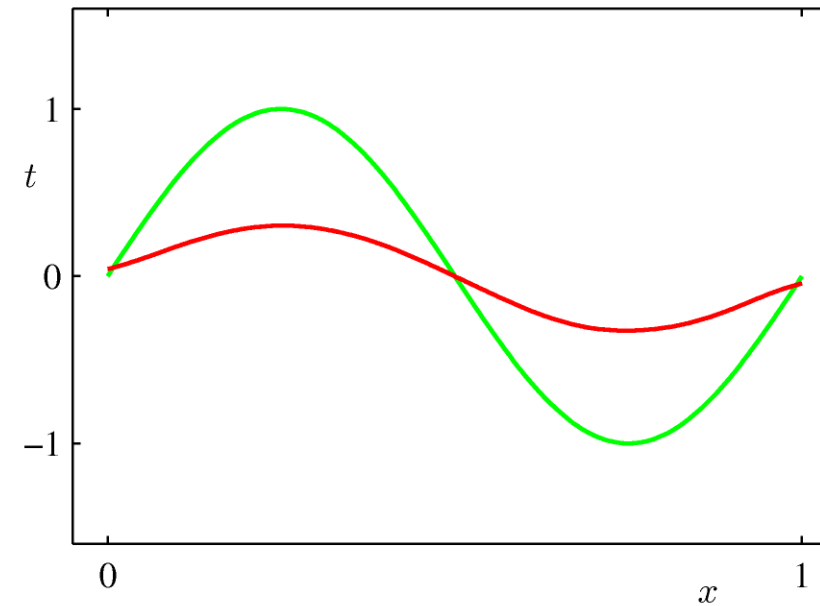
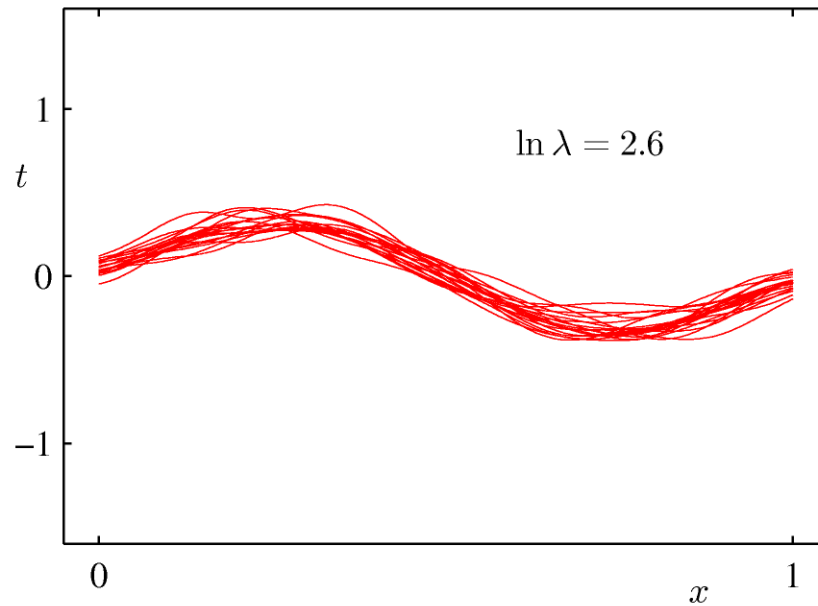
$$\textit{variance} = \int \mathbb{E}_{\mathcal{D}}[\{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})]\}^2] p(x) dx$$

$$\textit{noise} = \iint \{h(x) - t\}^2 p(x, t) dx dt$$



The Bias-Variance Decomposition₍₅₎

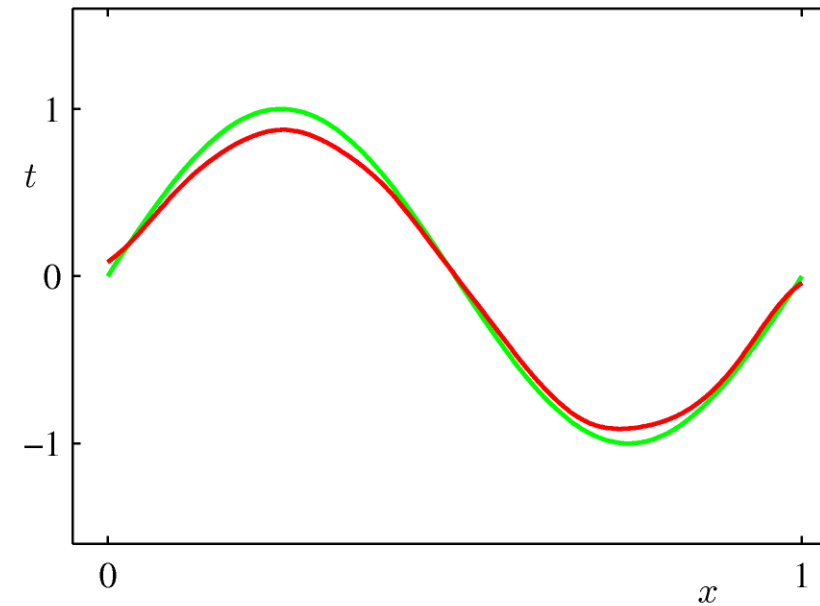
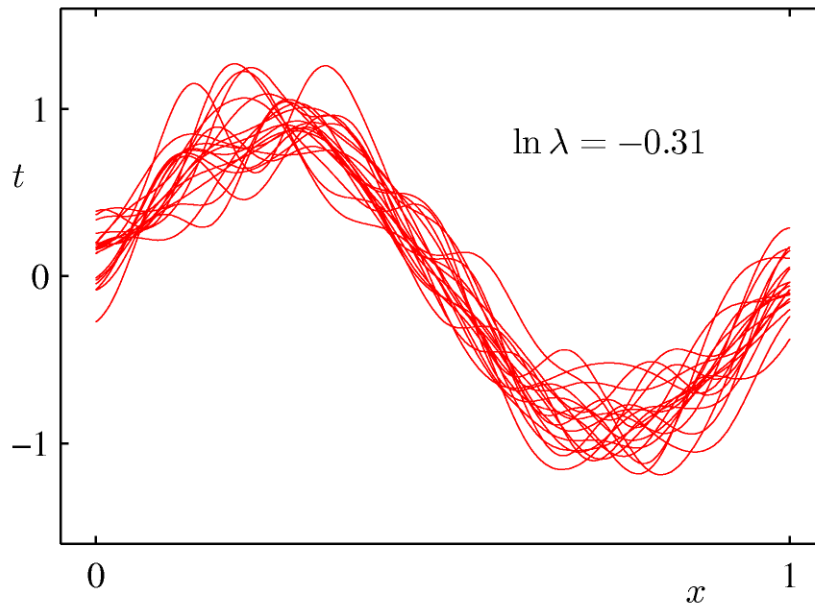
- Example: 25 data sets from the sinusoidal, varying the degree of regularization,





The Bias-Variance Decomposition₍₆₎

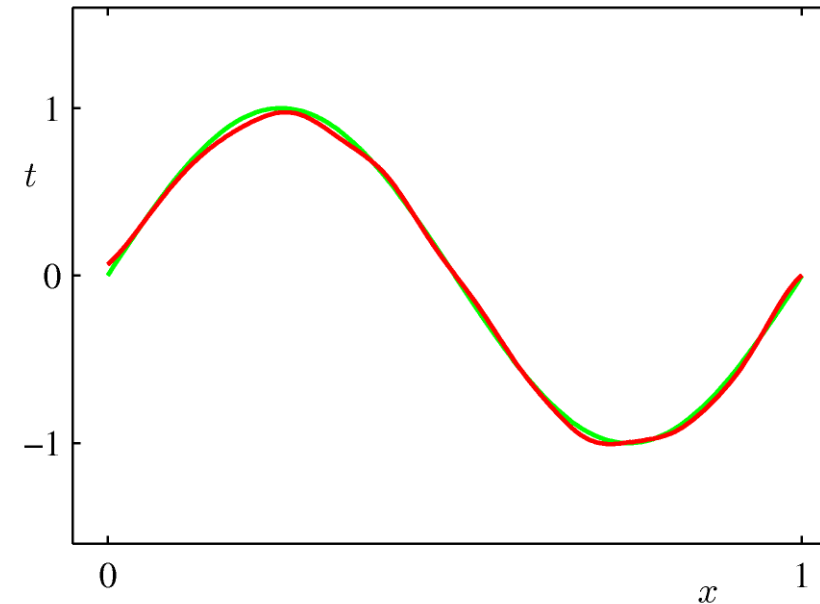
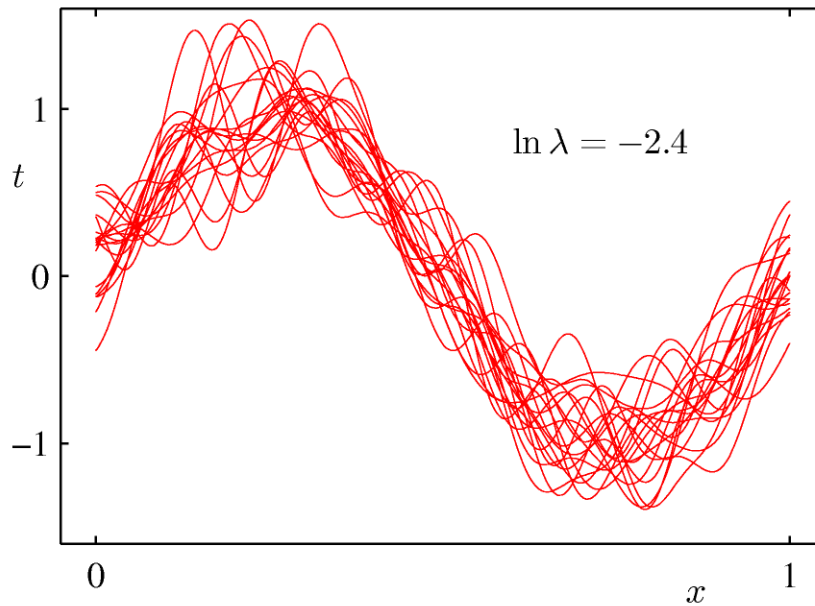
- Example: 25 data sets from the sinusoidal, varying the degree of regularization,





The Bias-Variance Decomposition₍₇₎

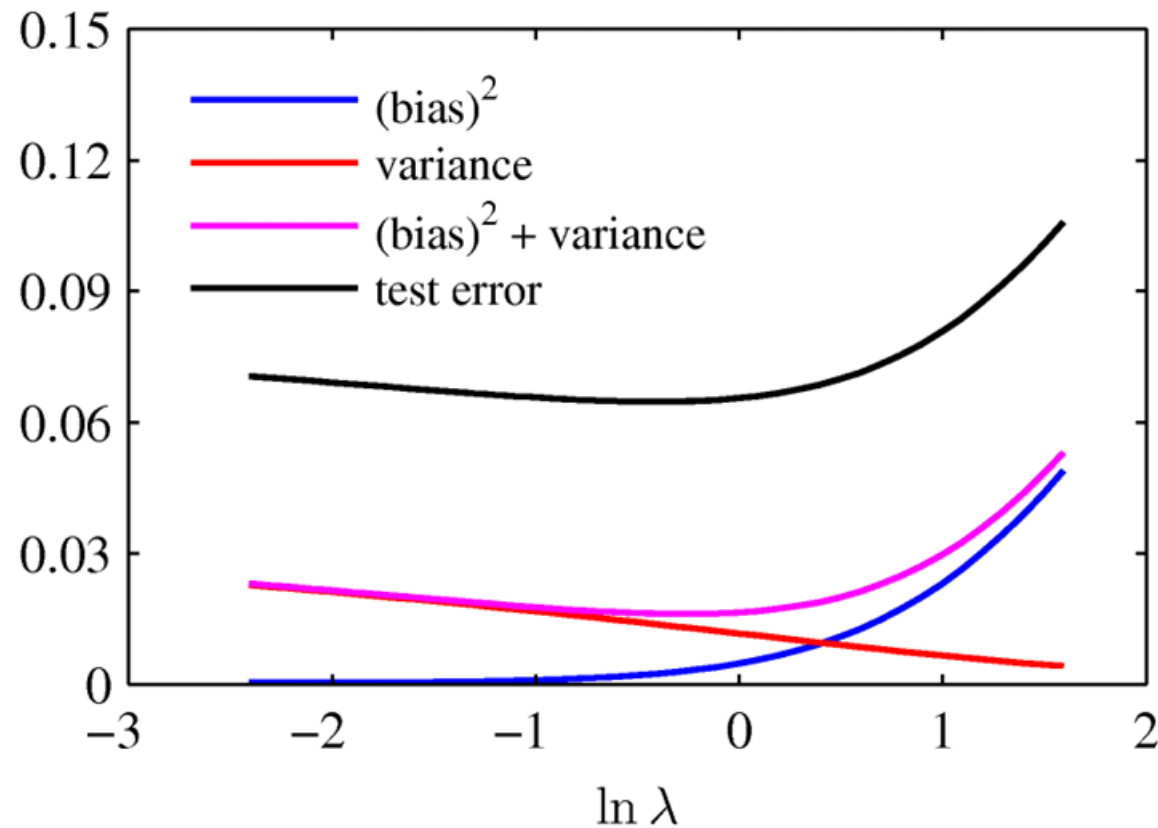
- Example: 25 data sets from the sinusoidal, varying the degree of regularization,





The Bias-Variance Trade-off

- From these plots, we note that an over-regularized model (large λ) will have a high bias, while an under-regularized model (small λ) will have a high variance.





Bayesian Linear Regression₍₁₎

- Define a conjugate prior over w

$$p(w) = \mathcal{N}(w|m_0, S_0)$$

Combining this with the likelihood function and using results for marginal and conditional Gaussian distributions, gives the posterior

where

$$p(w|t) = \mathcal{N}(w|m_N, S_N)$$

$$m_N = S_N(S_0^{-1}m_0 + \beta\Phi^T t)$$

$$S_N^{-1} = S_0^{-1} + \beta\Phi^T\Phi$$



Bayesian Linear Regression₍₂₎

- A common choice for the prior is

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

for which

$$m_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

Next we consider an example ...



Predictive Distribution₍₁₎

- Predict t for new values of x by integrating over w :

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

where

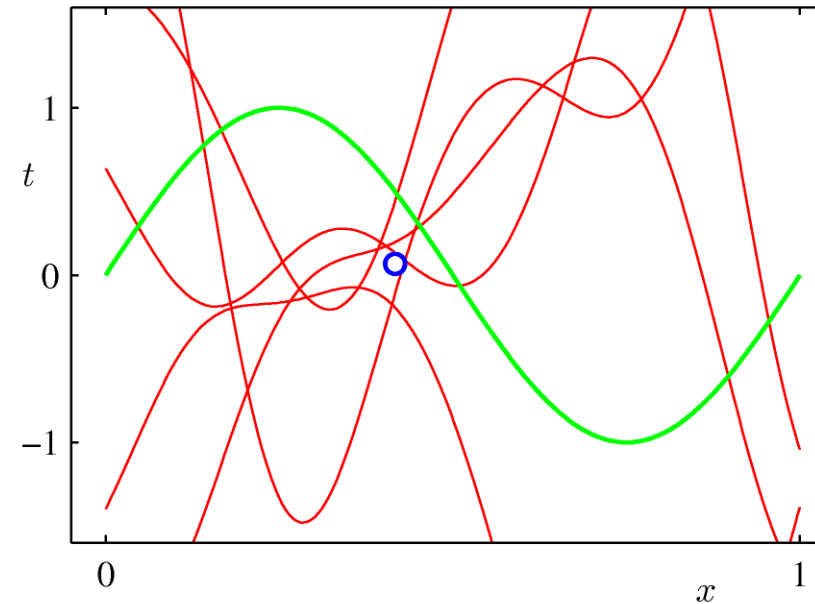
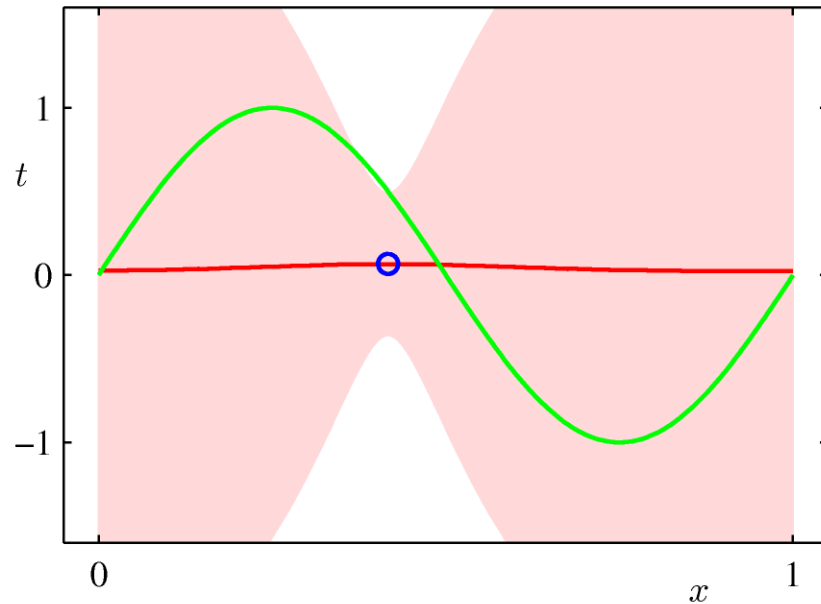
$$= \mathcal{N}(t|\mathbf{m}_N^T \phi(x), \sigma_N^2(x))$$

$$\sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N \phi(x)$$



Predictive Distribution₍₂₎

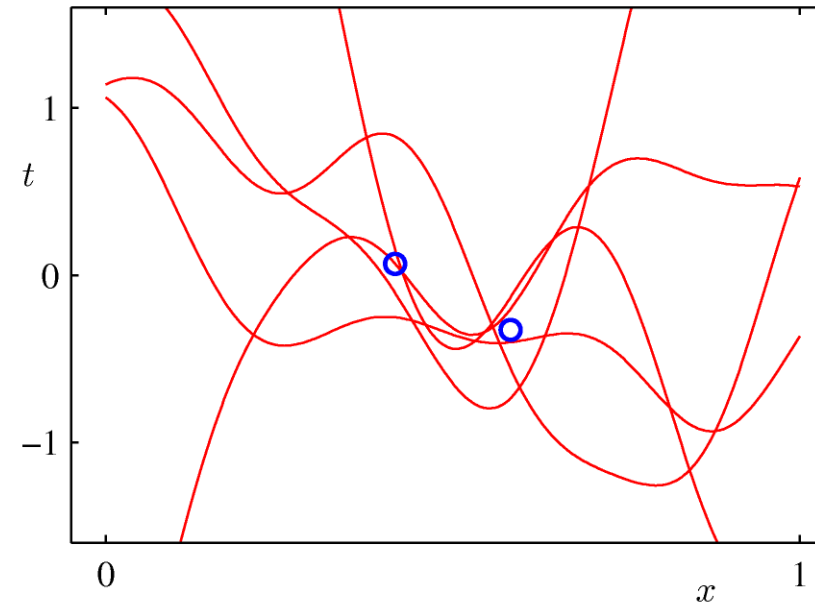
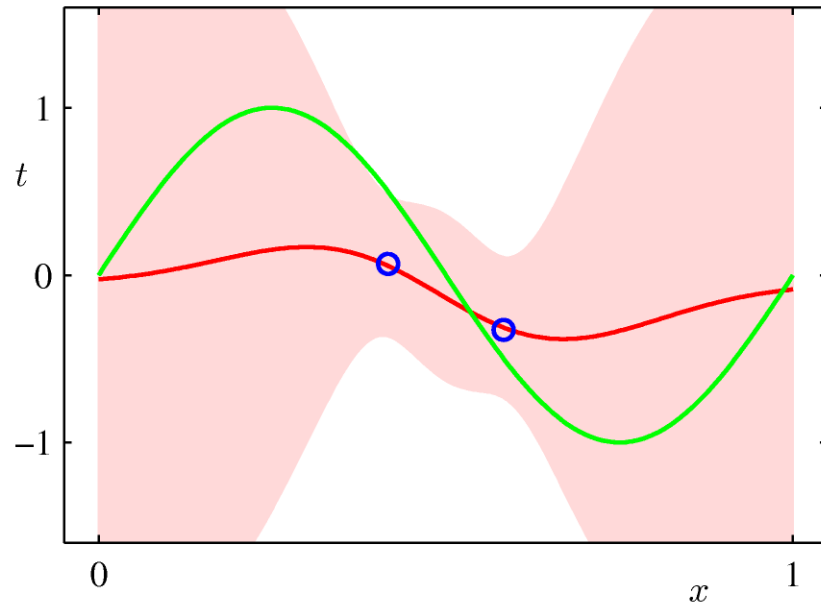
- Example: Sinusoidal data, 9 Gaussian basis functions, 1 data point





Predictive Distribution₍₃₎

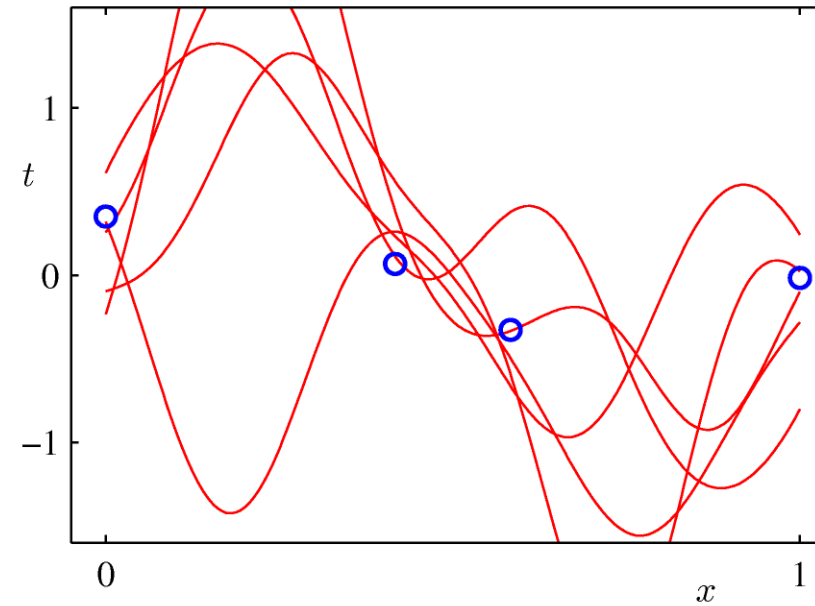
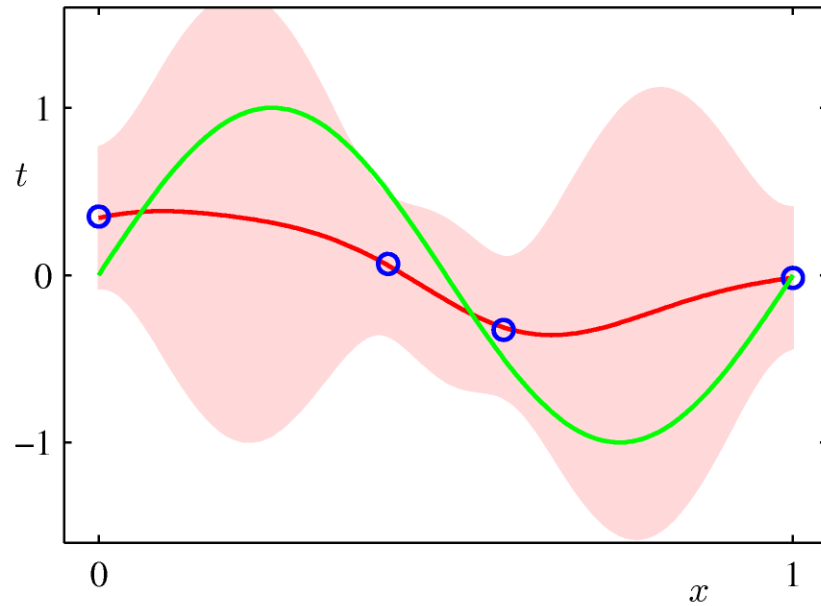
- Example: Sinusoidal data, 9 Gaussian basis functions, 2 data point





Predictive Distribution₍₄₎

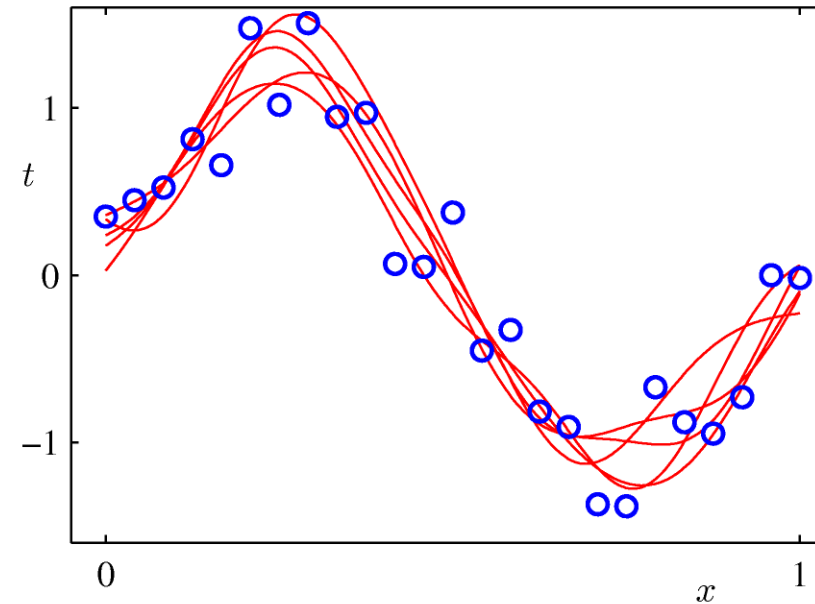
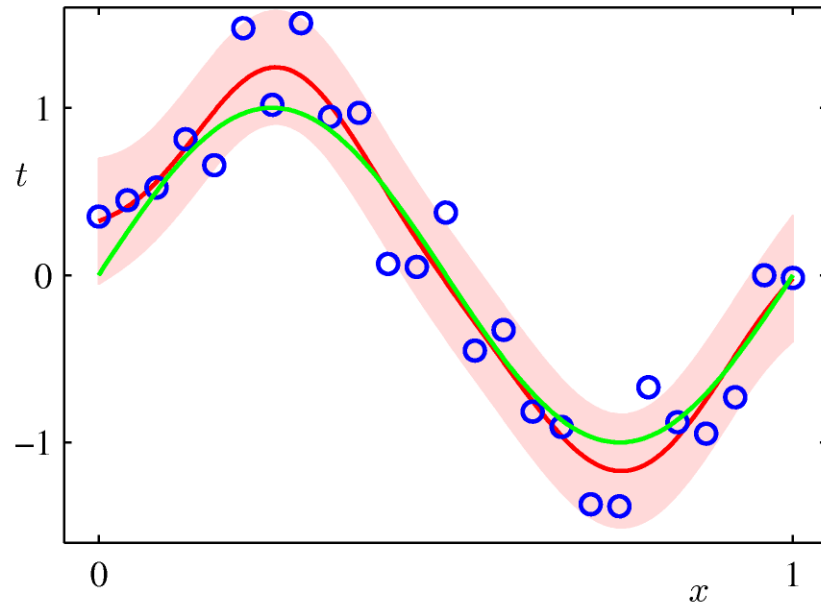
- Example: Sinusoidal data, 9 Gaussian basis functions, 4 data point





Predictive Distribution₍₄₎

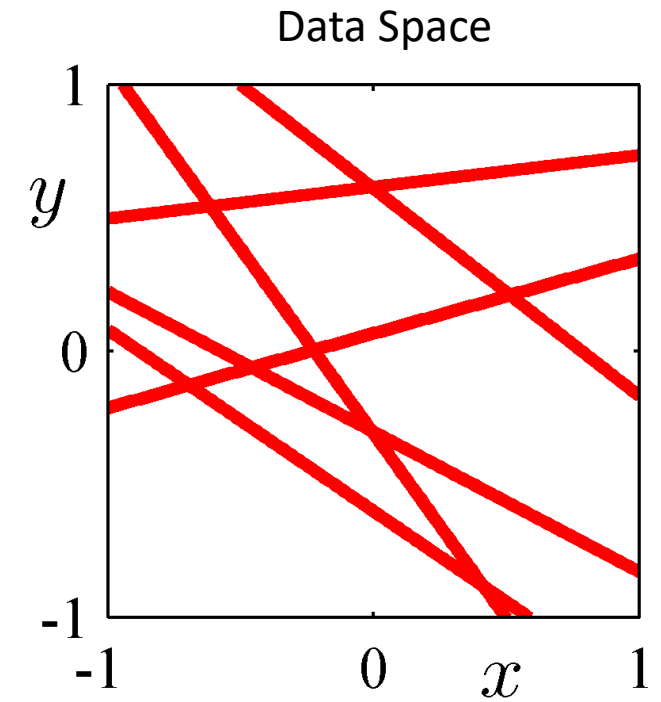
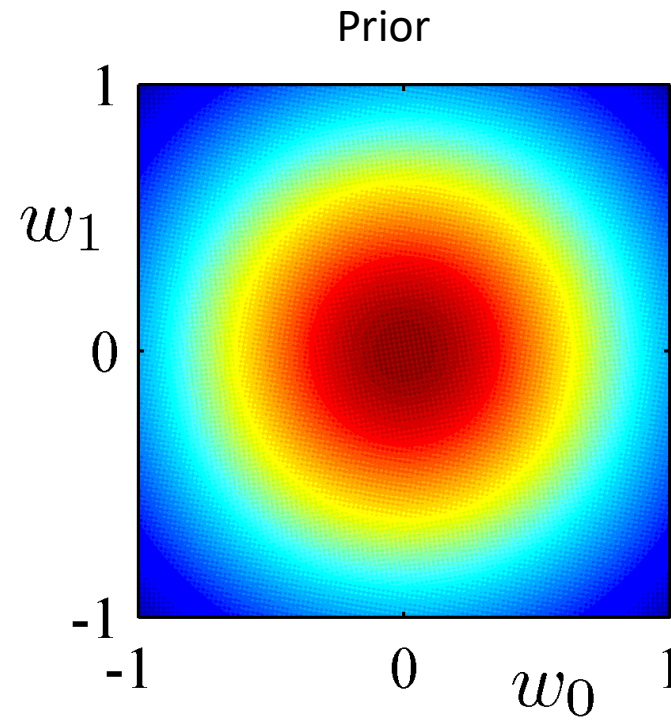
- Example: Sinusoidal data, 9 Gaussian basis functions, 25 data point





Bayesian Linear Regression₍₃₎

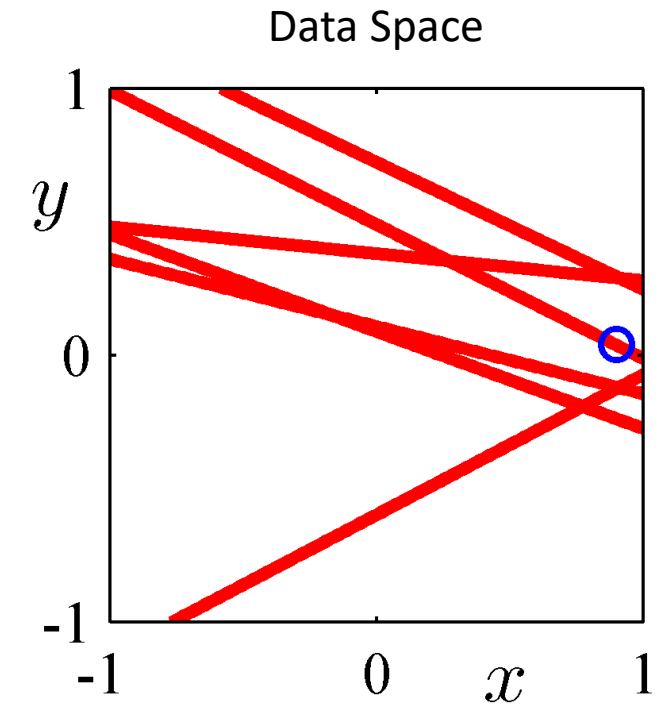
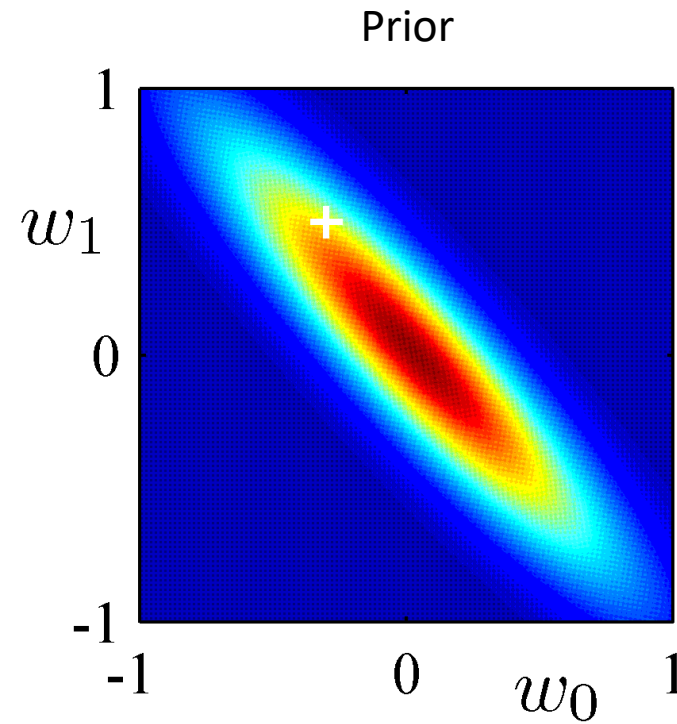
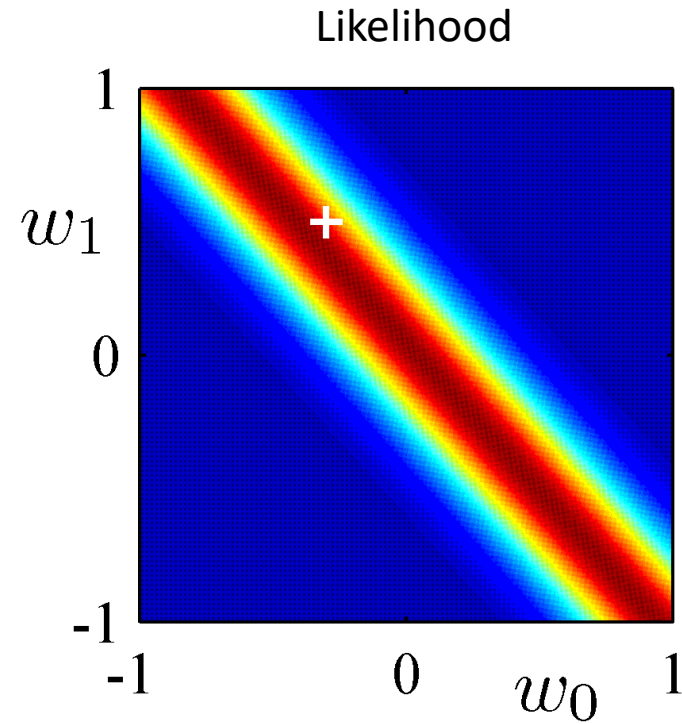
- 0 data points observed





Bayesian Linear Regression₍₄₎

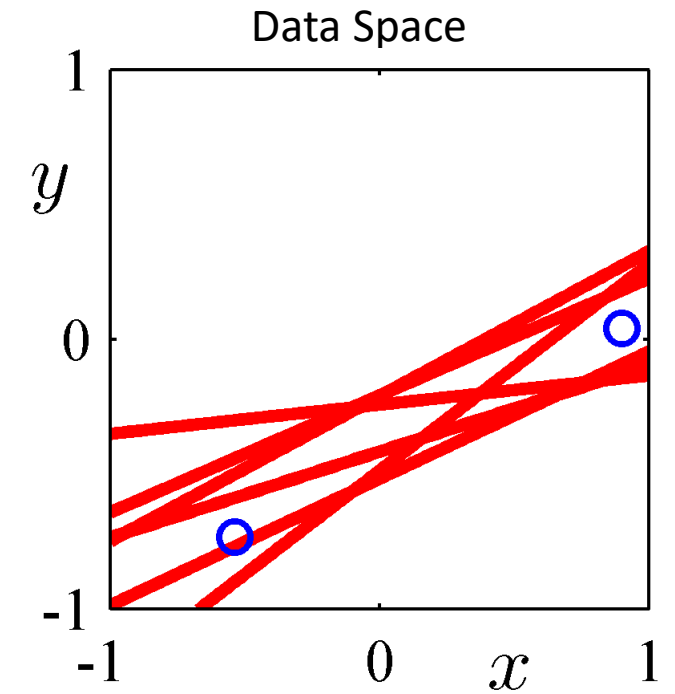
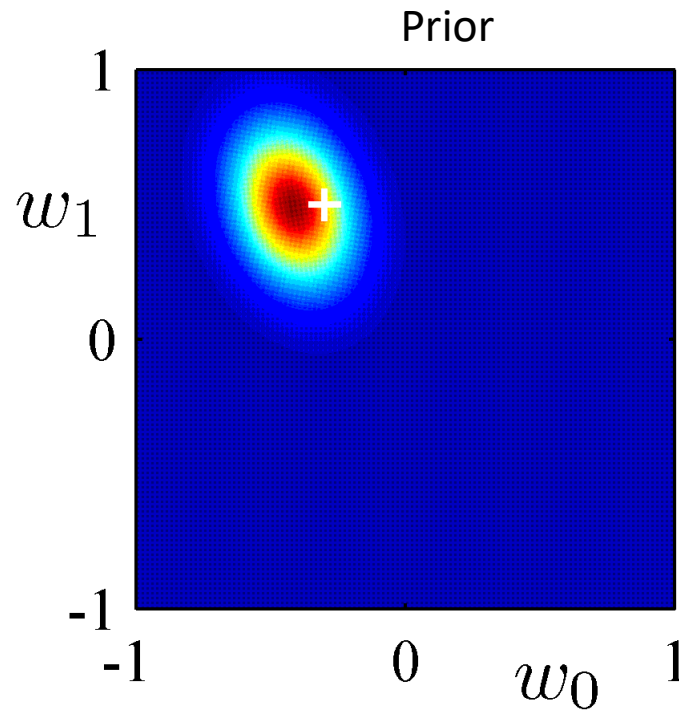
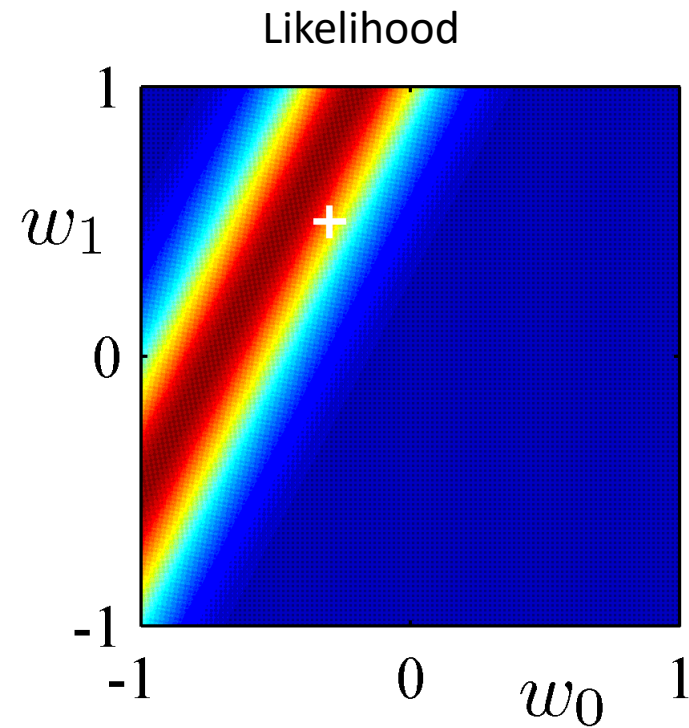
- 1 data points observed





Bayesian Linear Regression₍₅₎

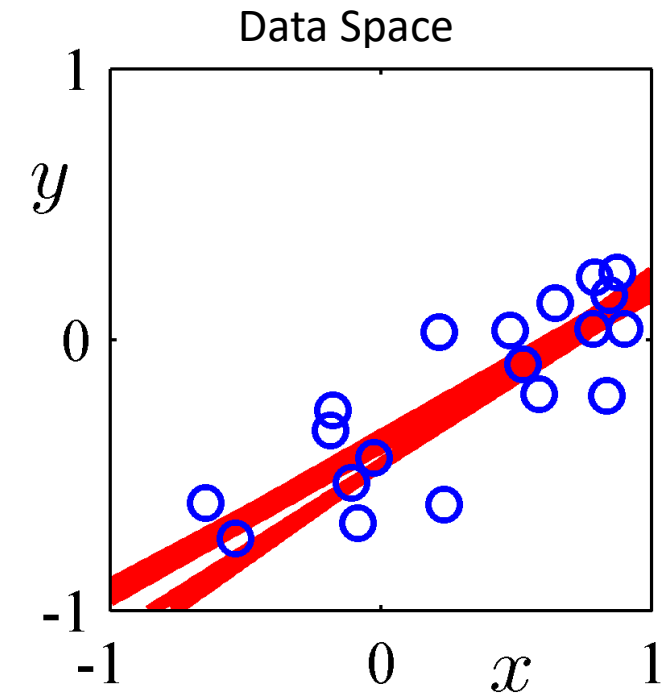
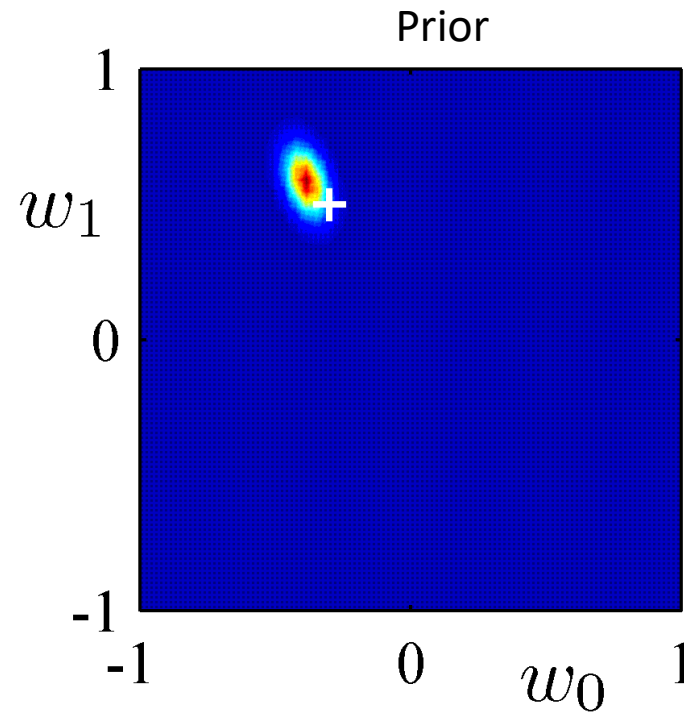
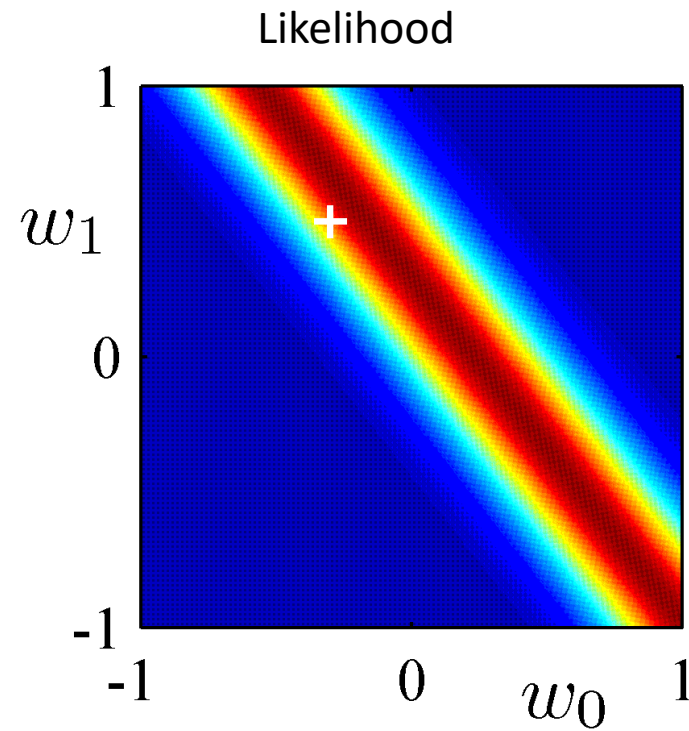
- 2 data points observed





Bayesian Linear Regression₍₆₎

- 20 data points observed





Equivalent Kernel ₍₁₎

- The predictive mean can be written

$$\begin{aligned}y(\mathbf{x}, \mathbf{m}_N) &= \mathbf{m}_N^T \phi(\mathbf{x}) = \beta \phi(\mathbf{x})^T \mathbf{S}_N \Phi^T \mathbf{t} \\&= \sum_{n=1}^N \underbrace{\beta \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}_n)}_{k(\mathbf{x}, \mathbf{x}_n)} t_n \\&= \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n.\end{aligned}$$

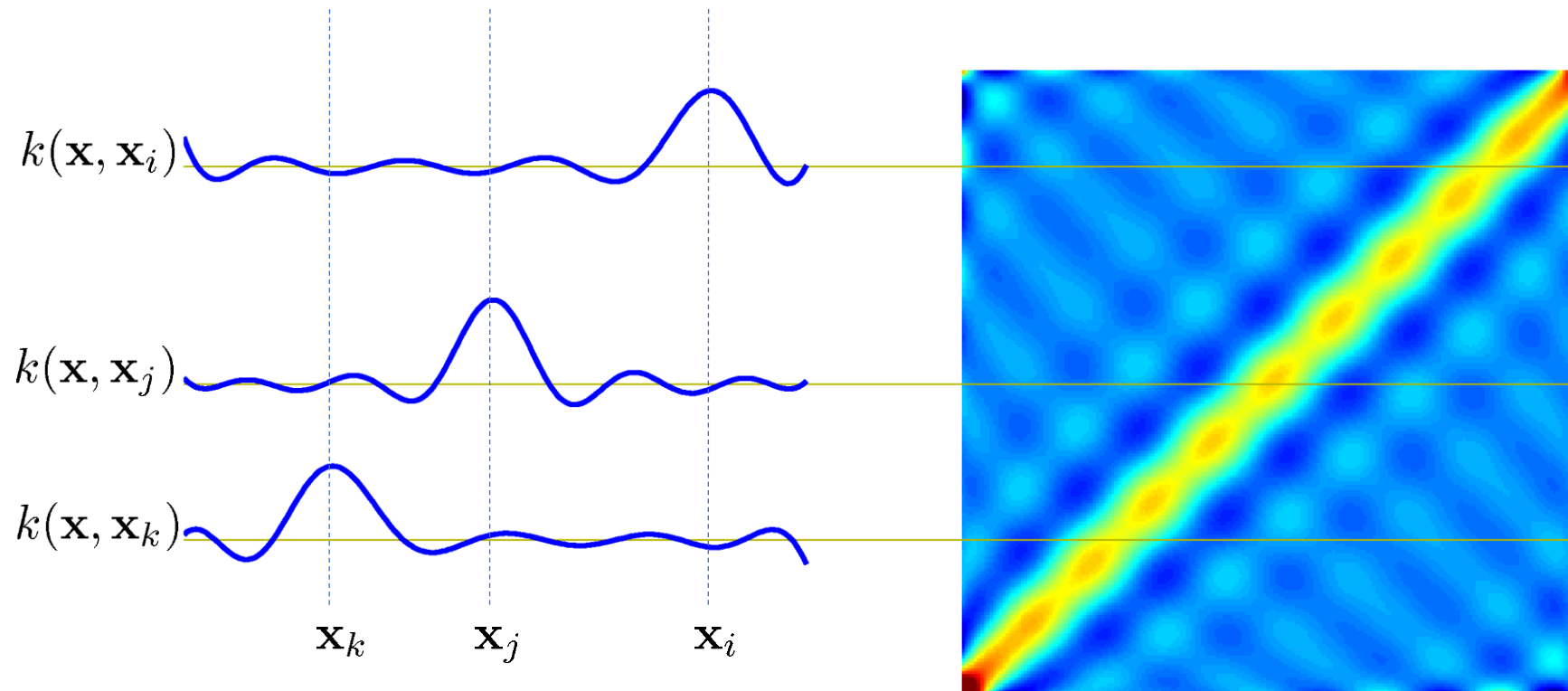
*Equivalent kernel or
smoother matrix.*

- This is a weighted sum of the training data target values, t_n .



Equivalent Kernel (2)

- Weight of t_n depends on distance between x and x_n ; nearby x_n carry more weight.





The Evidence Approximation ₍₁₎

- The fully Bayesian predictive distribution is given by

$$p(t|\mathbf{t}) = \iiint p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) p(\alpha, \beta|\mathbf{t}) d\mathbf{w} d\alpha d\beta$$

- but this integral is intractable. Approximate with

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

- where $(\hat{\alpha}, \hat{\beta})$ is the mode of $p(\alpha, \beta|\mathbf{t})$, which is assumed to be sharply peaked;
a.k.a. *empirical Bayes*, *type II* or *generalized maximum likelihood*, or *evidence approximation*.



The Evidence Approximation (2)

- From Bayes' theorem we have

$$p(\alpha, \beta | \mathbf{t}) \propto p(\mathbf{t} | \alpha, \beta) p(\alpha, \beta)$$

- and if we assume to $p(\alpha, \beta)$ be flat we see that

$$\begin{aligned} p(\alpha, \beta | \mathbf{t}) &\propto p(\mathbf{t} | \alpha, \beta) \\ &= \int p(\mathbf{t} | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) d\mathbf{w}. \end{aligned}$$

- General results for Gaussian integrals give

$$\ln p(\mathbf{t} | \alpha, \beta) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - E(\mathbf{m}_N) + \frac{1}{2} \ln |\mathbf{S}_N| - \frac{N}{2} \ln(2\pi).$$



Maximizing the Evidence Function ₍₁₎

- To maximise $p(t|\alpha, \beta)$ w.r.t. α and β , we define the eigenvector equation

$$\left(\beta \Phi^T \Phi\right) \mathbf{u}_i = \lambda_i \mathbf{u}_i.$$

- Thus

$$\mathbf{A} = \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

- has eigenvalues $\lambda_i + \alpha$.



- We can now differentiate $p(t|\alpha, \beta)$ w.r.t. α and β and set the results to zero, to get

$$\alpha = \frac{\gamma}{\mathbf{m}_N^T \mathbf{m}_N}$$

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

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

N.B. γ depends on both α and β

Any questions?

