



ISTA 421/521

Introduction to Machine Learning

Lecture 9: Maximum Likelihood Uncertainty 2

Clay Morrison

clayton@sista.arizona.edu

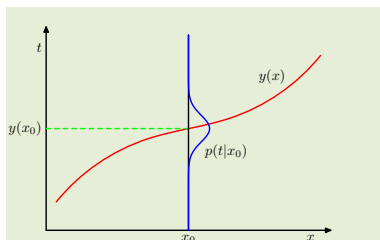
Gould-Simpson 819

Phone 621-6609

23 September 2014



Review

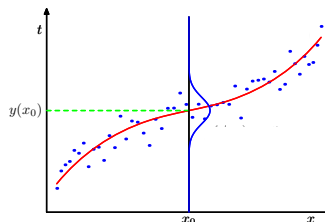


The generating process...

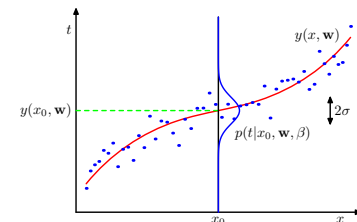
$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$L = p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \sigma^2)$$

$$= \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$



... generates data ...



... that we fit a model to

$$p(\hat{\mathbf{t}}|\mathbf{X}, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N p(\hat{t}_n|\mathbf{x}_n, \hat{\mathbf{w}}, \hat{\sigma}^2)$$

$$= \prod_{n=1}^N \mathcal{N}(\hat{\mathbf{w}}^\top \mathbf{x}_n, \hat{\sigma}^2)$$

prediction estimated parameters

Maximum Likelihood
Estimates of Params

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

The MLE is
unique

$$\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

Estimating Uncertainty
in Param Estimates via Expected Value

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = \mathbf{w}$$

The Fisher
Information $\mathcal{I} = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = -\left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top}\right)^{-1}$$

New Predictions:



Variability in Predictions

- We are predicting 2 values:

$$t_{new}, \sigma_{new}^2$$

$$t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new} \quad \text{Same solution as minimizing mean squared loss}$$

$$\begin{aligned} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\}^\top \mathbf{x}_{new} \\ &= \mathbf{w}^\top \mathbf{x}_{new} \end{aligned}$$

The **expected value** of our prediction is the new data attribute multiplied by the **true w**



Predicting the Variance of t_{new}

$$\sigma_{new}^2 = \text{var} \{t_{new}\} = \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}^2\} - (\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\})^2$$

Substitute $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$

$$\begin{aligned} \text{var} \{t_{new}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{(\hat{\mathbf{w}}^\top \mathbf{x}_{new})^2\} - (\mathbf{w}^\top \mathbf{x}_{new})^2 \\ &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{x}_{new}^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_{new}\} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}. \end{aligned}$$

Substitute $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$

$$\text{var} \{t_{new}\} = \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t} \mathbf{t}^\top\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}$$

On slide 22 of lec 8, in the derivation of the covariance of $\hat{\mathbf{w}}$, we identified $\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t} \mathbf{t}^\top\} = \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top + \sigma^2 \mathbf{I}$

$$\begin{aligned} \text{var} \{t_{new}\} &= \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I} + \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} \\ &= \sigma^2 \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} + \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} \\ &= \sigma^2 \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new}. \end{aligned}$$

Recall: $\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = - \left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right)^{-1}$

So, could be written

$$\sigma_{new}^2 = \mathbf{x}_{new}^\top \text{COV}\{\hat{\mathbf{w}}\} \mathbf{x}_{new}$$



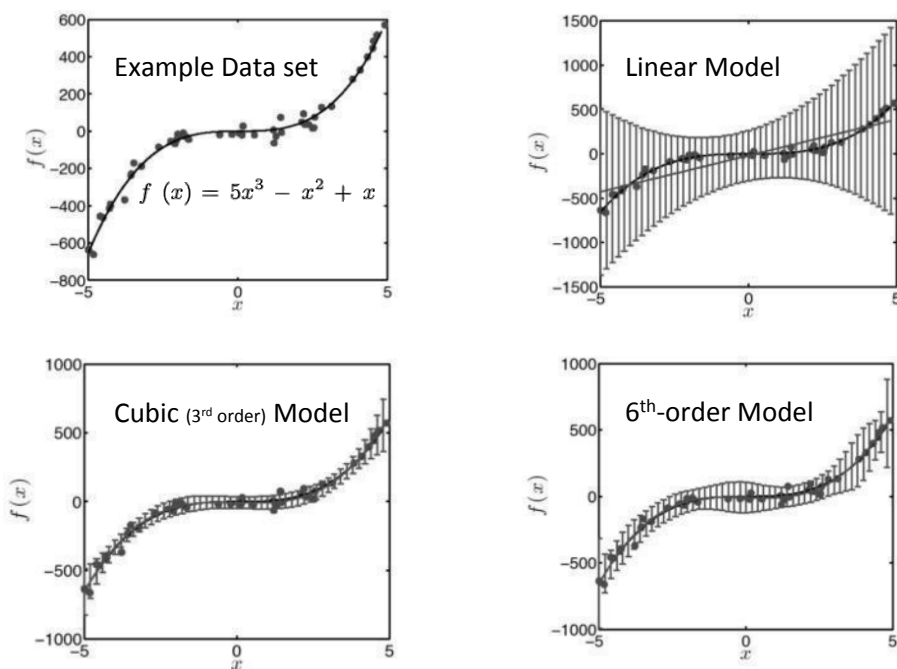
Prediction Summary

$$t_{\text{new}} = \mathbf{x}_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} = \mathbf{x}_{\text{new}}^T \hat{\mathbf{W}}$$

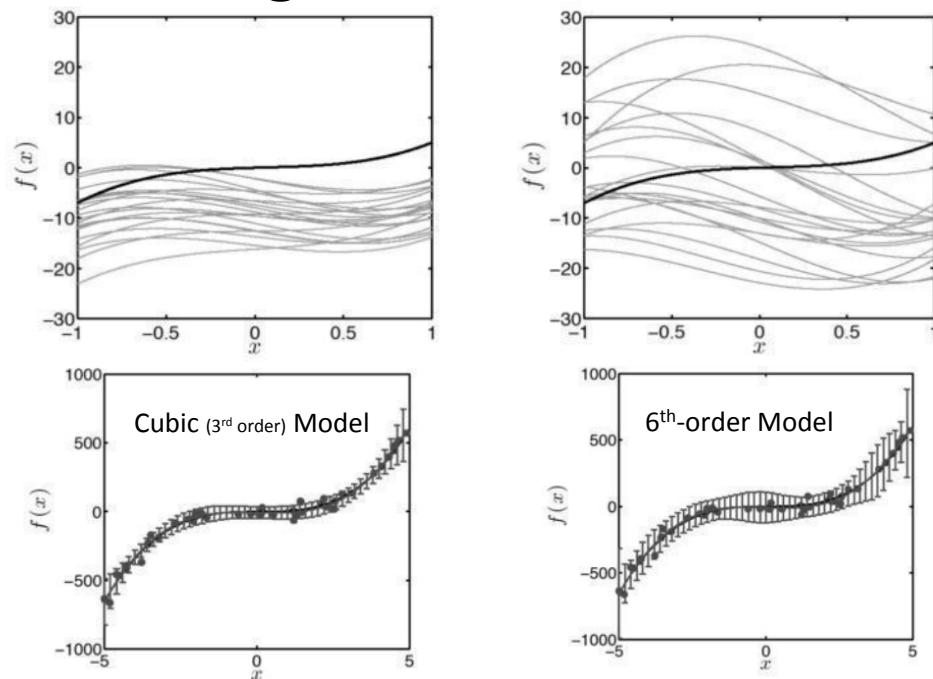
$$\sigma_{\text{new}}^2 = \sigma^2 \mathbf{x}_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{\text{new}}$$

We estimate this from the data: $\hat{\sigma}^2$

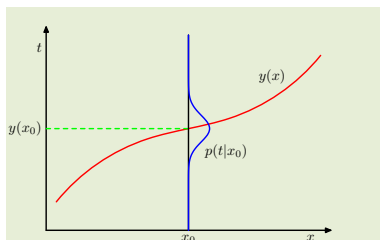
Plotting Predictive Error Bars



Plotting Predictive Error Bars



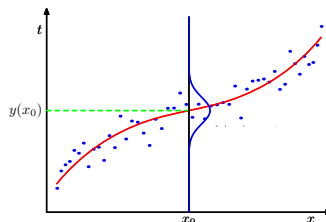
Review



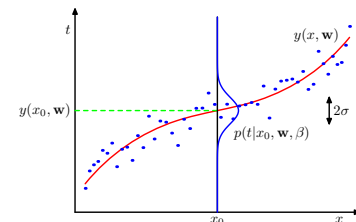
The generating process...

$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$L = p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$



... generates data ...



... that we fit a model to

$$p(\hat{\mathbf{t}}|\mathbf{X}, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N p(\hat{t}_n|\mathbf{x}_n, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N \mathcal{N}(\hat{\mathbf{w}}^\top \mathbf{x}_n, \hat{\sigma}^2)$$

prediction estimated parameters

Maximum Likelihood Estimates of Params

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

The MLE is unique

$$\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

Estimating Uncertainty in Param Estimates via Expected Value

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = \mathbf{w}$$

The Fisher Information $\mathcal{I} = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = -\left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right)^{-1}$$

New Predictions: $t_{\text{new}} = \hat{\mathbf{w}}^\top \mathbf{x}_{\text{new}}$

$$\sigma_{\text{new}}^2 = \sigma^2 \mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{new}}$$

$$\sigma_{\text{new}}^2 = \mathbf{x}_{\text{new}}^\top \text{cov}\{\hat{\mathbf{w}}\} \mathbf{x}_{\text{new}}$$

Quantifying the Uncertainty in our Estimate of $\hat{\sigma}^2$

$$\widehat{\sigma^2} = \frac{1}{N}(\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

$$\begin{aligned} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \widehat{\sigma^2} \right\} &= \frac{1}{N} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}} \right\} \\ &= \frac{1}{N} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \right\} \\ &= \frac{1}{N} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{t}^\top \mathbf{t} \right\} - \frac{1}{N} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} \right\} \end{aligned}$$

We have seen the expectation of $\mathbf{t} \mathbf{t}^\top$... but not $\mathbf{t}^\top \mathbf{t}$

$$\begin{aligned} \mathbf{t} &\sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \mathbf{E}_{p(\mathbf{t})} \left\{ \mathbf{t}^\top \mathbf{A} \mathbf{t} \right\} &= \text{Tr}(\mathbf{A} \boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu} \end{aligned}$$



Quantifying the Uncertainty in our Estimate of $\hat{\sigma}^2$

- We will need the matrix **trace**

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1D} \\ A_{21} & A_{22} & \cdots & A_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ A_{D1} & A_{D2} & \cdots & A_{DD} \end{bmatrix}$$

$$\text{Tr}(\mathbf{A}) = \sum_{d=1}^D A_{dd}$$

Other Properties/Identities

$$\text{Tr}(\mathbf{I}_D) = \sum_{d=1}^D 1 = D$$

$$\text{Tr}(\mathbf{A}\mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{A})$$

$$\text{Tr}(a) = a$$

$$\text{Tr}(\mathbf{w}^\top \mathbf{w}) = \mathbf{w}^\top \mathbf{w}$$

- The trace of a matrix is the sum of the (complex) eigenvalues, and is invariant with respect to change of basis.
- Geometrically: The trace can be interpreted as the infinitesimal change in volume (as the derivative of the determinant)



Quantifying the Uncertainty in our Estimate of $\hat{\sigma}^2$

$$\mathbf{t} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

$$\mathbb{E}_{p(\mathbf{t})} \{ \mathbf{t}^T \mathbf{A} \mathbf{t} \} = \text{Tr}(\mathbf{A} \Sigma) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$$

$$\text{Tr}(\mathbf{A}) = \sum_{d=1}^D A_{dd}$$

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\sigma}^2 \} = \frac{1}{N} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \mathbf{t}^T \mathbf{t} \} - \frac{1}{N} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \mathbf{t}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \}$$

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\sigma}^2 \} = \frac{1}{N} \left(\text{Tr}(\sigma^2 \mathbf{I}_N) + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \right) - \frac{1}{N} \left(\text{Tr}(\sigma^2 \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) + \mathbf{w}^T \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w} \right)$$

$$\text{Tr}(\sigma^2 \mathbf{A}) = \sigma^2 \text{Tr}(\mathbf{A}) \text{ and } \text{Tr}(\mathbf{I}_N) = N$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\sigma}^2 \} &= \sigma^2 + \frac{1}{N} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \frac{\sigma^2}{N} \text{Tr}(\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) - \frac{1}{N} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \\ &= \sigma^2 - \frac{\sigma^2}{N} \text{Tr}(\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \\ &= \sigma^2 \left(1 - \frac{1}{N} \text{Tr}(\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \right). \end{aligned}$$

$$\text{Tr}(\mathbf{A} \mathbf{B}) = \text{Tr}(\mathbf{B} \mathbf{A})$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\sigma}^2 \} &= \sigma^2 \left(1 - \frac{1}{N} \text{Tr}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}) \right) \\ &= \sigma^2 \left(1 - \frac{1}{N} \text{Tr}(\mathbf{I}_D) \right) \\ &= \sigma^2 \left(1 - \frac{D}{N} \right), \end{aligned}$$

D is the number of attributes (the number of columns in \mathbf{X})



11

Quantifying the Uncertainty in our Estimate of $\hat{\sigma}^2$

$$\mathbf{t} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

$$\mathbb{E}_{p(\mathbf{t})} \{ \mathbf{t}^T \mathbf{A} \mathbf{t} \} = \text{Tr}(\mathbf{A} \Sigma) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$$

$$\text{Tr}(\mathbf{A}) = \sum_{d=1}^D A_{dd}$$

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\sigma}^2 \} = \frac{1}{N} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \mathbf{t}^T \mathbf{t} \} - \frac{1}{N} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \mathbf{t}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \}$$

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\sigma}^2 \} = \sigma^2 \left(1 - \frac{D}{N} \right)$$

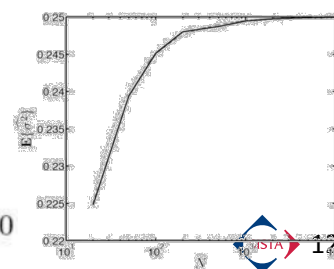
When $D < N$ (that is, the number of attributes we measure for each data point is *smaller* than the number of data points), then our estimates of the variance will, on average, be lower than the true variance.

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\sigma}^2 \} < \sigma^2$$

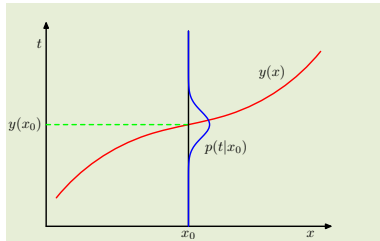
Unlike the estimate for $\hat{\mathbf{w}}$, the MLE for $\hat{\sigma}^2$ is *biased*.

$$D = 2 \text{ and } N = 20$$

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\sigma}^2 \} = \sigma^2 \left(1 - \frac{D}{N} \right) = 0.25 \left(1 - \frac{2}{20} \right) = 0.2250$$



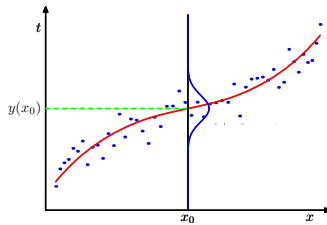
Review



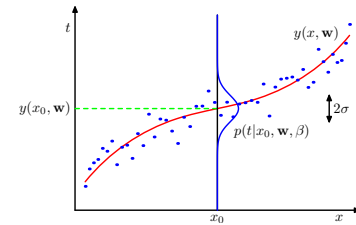
The generating process...

$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$L = p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$



... generates data ...



... that we fit a model to

$$p(\hat{\mathbf{t}}|\mathbf{X}, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N p(\hat{t}_n|\mathbf{x}_n, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N \mathcal{N}(\hat{\mathbf{w}}^\top \mathbf{x}_n, \hat{\sigma}^2)$$

prediction estimated parameters

Maximum Likelihood Estimates of Params

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

The MLE is unique

$$\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

Estimating Uncertainty in Param Estimates via Expected Value

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = \mathbf{w}$$

The Fisher Information

$$\mathcal{I} = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = -\left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top}\right)^{-1}$$

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\sigma}^2\} = \sigma^2 \left(1 - \frac{D}{N}\right)$$

New Predictions: $t_{\text{new}} = \hat{\mathbf{w}}^\top \mathbf{x}_{\text{new}}$

$$\sigma_{\text{new}}^2 = \sigma^2 \mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{new}}$$

$$\sigma_{\text{new}}^2 = \mathbf{x}_{\text{new}}^\top \text{cov}\{\hat{\mathbf{w}}\} \mathbf{x}_{\text{new}}$$

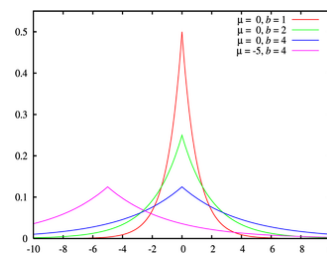
MLE Beyond the Gaussian

- Other useful probability distributions have analytic MLE solutions.

- The MLE of **Laplace** i.i.d. variables.

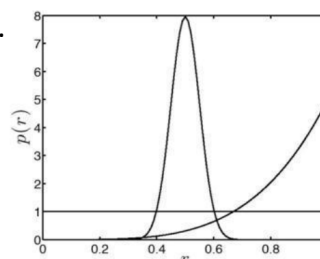
$$f(x|\mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$

$$= \frac{1}{2b} \begin{cases} \exp\left(-\frac{\mu - x}{b}\right) & \text{if } x < \mu \\ \exp\left(-\frac{x - \mu}{b}\right) & \text{if } x \geq \mu \end{cases}$$



- μ = "location" parameter, b = "diversity"
- Another univariate case is the **beta** distribution – also has a closed form, with lots of flexibility in the shape/location.

$$p(r) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} r^{\alpha-1} (1-r)^{\beta-1}$$

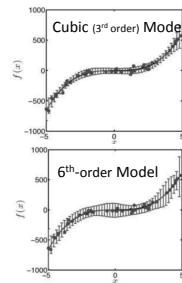


MLE and Model Selection

$$\log L = -\frac{1}{N} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \hat{\mathbf{w}})^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$



MLE Prefers Complex Models

$$\log L = -\frac{1}{N} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \hat{\mathbf{w}})^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

Plug in $\hat{\sigma}^2$ to the log likelihood:

$$\begin{aligned} \log L &= -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \hat{\sigma}^2 - \frac{1}{2\hat{\sigma}^2} N \hat{\sigma}^2 \\ &= -\frac{N}{2} (1 - \log 2\pi) - \frac{N}{2} \log \hat{\sigma}^2 \end{aligned}$$

Making $\hat{\sigma}^2$ smaller makes $\log L$ larger.

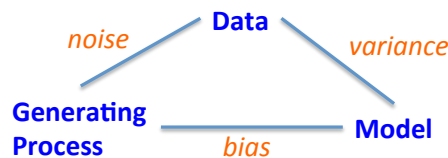
Bad news: Increasing the model complexity will *decrease* the variance!

Bottom line: Unfortunately, we can't use MLE to do model selection

But, with a particular model, MLE will choose the parameters that make the data have the highest overall likelihood under the model.

Bias-Variance Tradeoff

- **Bias:** the systematic mismatch between our model and the process that generated the data.
 - Too simple a model == too high a bias (underfitting)
 - Too complex a model (too many degrees of freedom) == too low a bias (overfitting)
- **Variance:** Squared error between model and data
- Imagine we had the *true* distribution that generated the data; we could compute and compare the expected value of the squared error between estimated parameter values and the true values.
- We would like this value to be as small as possible.



The Bias-Variance Decomposition (1)

$$\mathbb{E}_{p(x)} \{f(x)\} = \int f(x)p(x) dx$$

Sorry to switch Expectation notation on you...
 ← to the left is from our class book
 The notation below and the next couple of slides comes from Bishop (2007)

- Recall the *expected squared loss*,

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

- where

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

- The second term of $\mathbb{E}[L]$ corresponds to the noise inherent in the random variable t .
- What about the first term?

The Bias-Variance Decomposition (2)

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$
$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt.$$

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

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

The Bias-Variance Decomposition (3)

- Taking the expectation over \mathcal{D} yields

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

The Bias-Variance Decomposition (4)

- Thus we can write

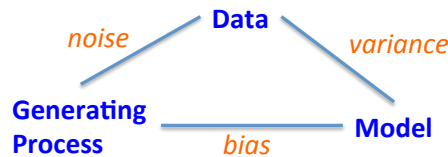
$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

- where

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

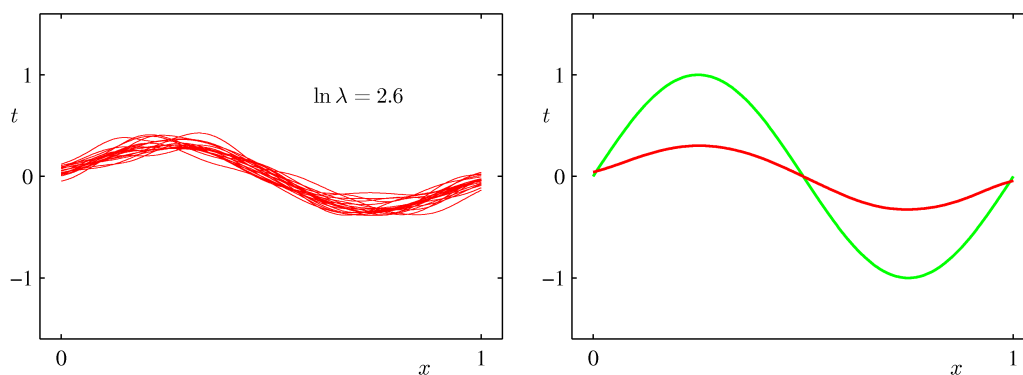
$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$



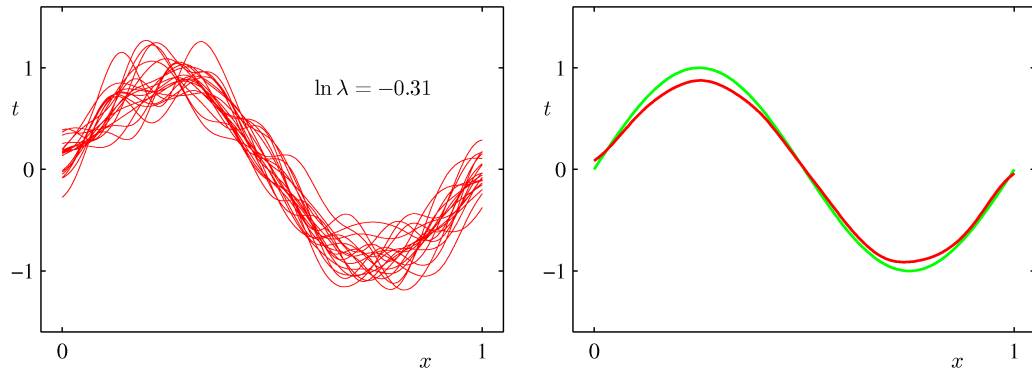
The Bias-Variance Decomposition (5)

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



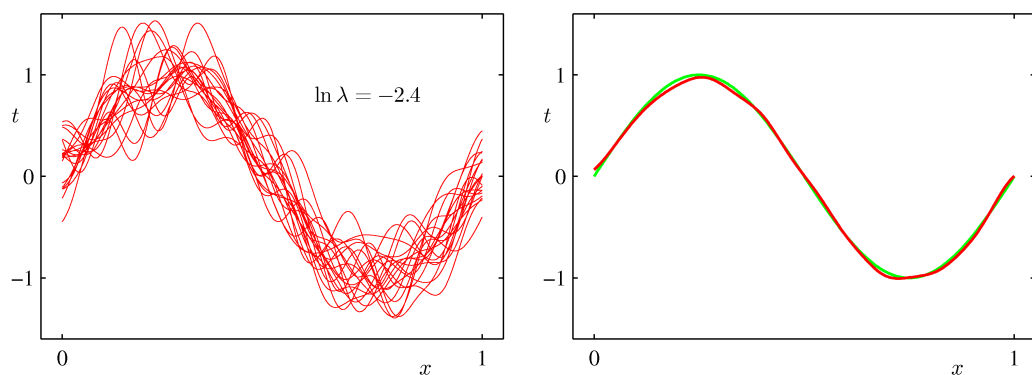
The Bias-Variance Decomposition (6)

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



The Bias-Variance Decomposition (7)

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



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.

