

Machine Learning Bias–Variance Tradeoff

Devdatt Dubhashi
dubhashi@chalmers.se

Dept. of Computer Science and Engg.
Chalmers University

March 26, 2018



- Introduction
- D. Dubhashi
- Confidence in parameter estimates
- Predictions
- Prediction
- Likelihood for model selection
- Summary

Optimum parameters

- ▶ We have point estimates of our parameters.
- ▶ How confident should we be in them?
 - ▶ If we changed them a little bit, would the model still be good?



- Introduction
- D. Dubhashi
- Confidence in parameter estimates
- Predictions
- Prediction
- Likelihood for model selection
- Summary

Confidence in parameter estimates

- ▶ Imagine there are **true** parameters, \mathbf{w} and σ^2 .



- Introduction
- D. Dubhashi
- Confidence in parameter estimates
- Predictions
- Prediction
- Likelihood for model selection
- Summary

Confidence in parameter estimates

- ▶ Imagine there are **true** parameters, \mathbf{w} and σ^2 .
- ▶ How good are our estimates $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$?
 - ▶ Are they correct (on average)?
 - ▶ If we could keep adding data, would we converge on the true value?



- Introduction
- D. Dubhashi
- Confidence in parameter estimates
- Predictions
- Prediction
- Likelihood for model selection
- Summary

Confidence in parameter estimates

- ▶ Imagine there are **true** parameters, \mathbf{w} and σ^2 .
- ▶ How good our our estimates $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$?
 - ▶ Are they correct (on average)?
 - ▶ If we could keep adding data, would we converge on the true value?
- ▶ How confident should we be in our estimates?
 - ▶ Could we change parameters a little bit and still have a good model?



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

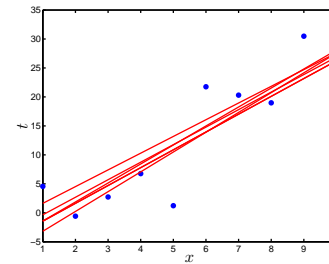
Prediction

Likelihood for model selection

Summary

Confidence in parameter estimates

- ▶ Imagine there are **true** parameters, \mathbf{w} and σ^2 .
- ▶ How good our our estimates $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$?
 - ▶ Are they correct (on average)?
 - ▶ If we could keep adding data, would we converge on the true value?
- ▶ How confident should we be in our estimates?
 - ▶ Could we change parameters a little bit and still have a good model?



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

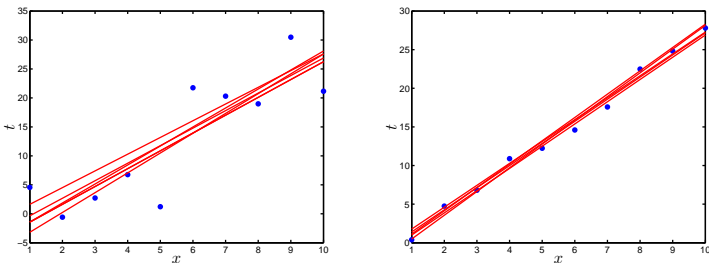
Prediction

Likelihood for model selection

Summary

Confidence in parameter estimates

- ▶ Imagine there are **true** parameters, \mathbf{w} and σ^2 .
- ▶ How good our our estimates $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$?
 - ▶ Are they correct (on average)?
 - ▶ If we could keep adding data, would we converge on the true value?
- ▶ How confident should we be in our estimates?
 - ▶ Could we change parameters a little bit and still have a good model?



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Predictions

- ▶ Our aim is to make predictions (e.g. London 2012)



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

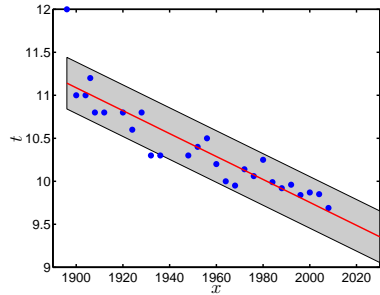
Prediction

Likelihood for model selection

Summary

Predictions

- ▶ Our aim is to make predictions (e.g. London 2012)
- ▶ The noise in our data tells us that we can't predict exactly.



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

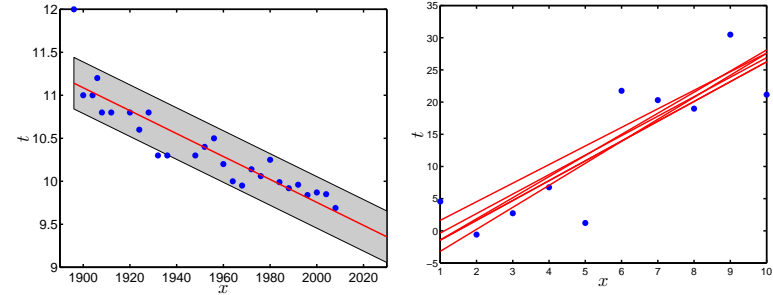
Prediction

Likelihood for model selection

Summary

Predictions

- ▶ Our aim is to make predictions (e.g. London 2012)
- ▶ The noise in our data tells us that we can't predict exactly.
- ▶ The uncertainty in the parameters $\text{cov}\{\hat{\mathbf{w}}\}$ should make them even less certain.



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

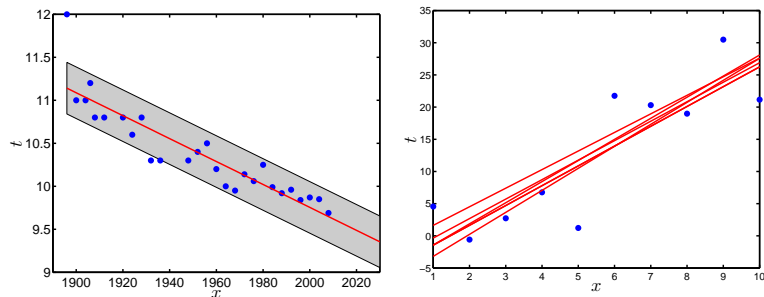
Prediction

Likelihood for model selection

Summary

Predictions

- ▶ Our aim is to make predictions (e.g. London 2012)
- ▶ The noise in our data tells us that we can't predict exactly.
- ▶ The uncertainty in the parameters $\text{cov}\{\hat{\mathbf{w}}\}$ should make them even less certain.



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Predictions

- ▶ Our model is defined as:
$$t = \mathbf{w}^T \mathbf{x} + \epsilon$$
- ▶ Given our estimate of the parameters, $\hat{\mathbf{w}}$ and a new input, \mathbf{x}_{new} , if we had to predict a single value:

$$t_{\text{new}} = \hat{\mathbf{w}}^T \mathbf{x}_{\text{new}}$$

- ▶ Is this sensible?



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

- ▶ We'll also (hopefully) clear up some loose ends.

Predictions

- ▶ Our model is defined as:

$$t = \mathbf{w}^T \mathbf{x} + \epsilon$$

- ▶ Given our estimate of the parameters, $\hat{\mathbf{w}}$ and a new input, \mathbf{x}_{new} , if we had to predict a single value:

$$t_{\text{new}} = \hat{\mathbf{w}}^T \mathbf{x}_{\text{new}}$$

- ▶ Is this sensible? What is $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{t_{\text{new}}\}$?

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

- ▶ which is a good thing!



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Predictions

- ▶ What about $\text{var}\{t_{\text{new}}\}$?

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



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Predictions

- ▶ What about $\text{var}\{t_{\text{new}}\}$?

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

$$= \mathbf{E} \{ (\hat{\mathbf{w}}^T \mathbf{x}_{\text{new}})^2 \} - (\mathbf{w}^T \mathbf{x}_{\text{new}})^2$$

$$= \mathbf{x}_{\text{new}}^T \mathbf{E} \{ \hat{\mathbf{w}} \hat{\mathbf{w}}^T \} \mathbf{x}_{\text{new}} - \mathbf{x}_{\text{new}}^T \mathbf{w} \mathbf{w}^T \mathbf{x}_{\text{new}}$$

= :

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



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Prediction and variance

$$t_{\text{new}} = \hat{\mathbf{w}}^T \mathbf{x}_{\text{new}}$$
$$\text{var}\{t_{\text{new}}\} = \sigma^2 \mathbf{x}_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{\text{new}}$$



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Prediction and variance

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

- ▶ Recall the expression for the covariance of the parameter estimate:

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Prediction and variance

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

- ▶ Recall the expression for the covariance of the parameter estimate:

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

- ▶ Appears in the variance of the prediction:

$$\text{var}\{t_{\text{new}}\} = \mathbf{x}_{\text{new}}^T \text{cov}\{\hat{\mathbf{w}}\} \mathbf{x}_{\text{new}}$$



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Prediction and variance

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

- ▶ Recall the expression for the covariance of the parameter estimate:

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

- ▶ Appears in the variance of the prediction:

$$\text{var}\{t_{\text{new}}\} = \mathbf{x}_{\text{new}}^T \text{cov}\{\hat{\mathbf{w}}\} \mathbf{x}_{\text{new}}$$

- ▶ If the variance in the parameters is high, so is the variance in the predictions.



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Example

Data sampled from a 3rd order polynomial function:

$$t = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \epsilon$$



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

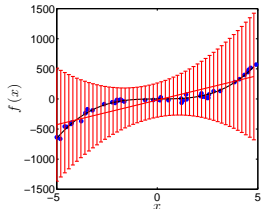
Likelihood for model selection

Summary

Example

Data sampled from a 3rd order polynomial function:

$$t = w_0 + w_1x + w_2x^2 + w_3x^3 + \epsilon$$



Linear

Plots show $t_{\text{new}} \pm \text{var}\{t_{\text{new}}\}$. (Black line is truth).



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

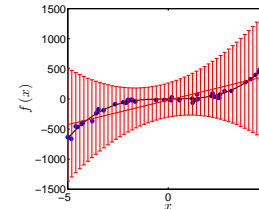
Likelihood for model selection

Summary

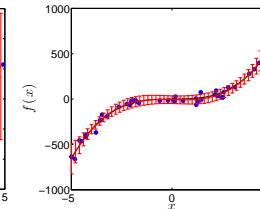
Example

Data sampled from a 3rd order polynomial function:

$$t = w_0 + w_1x + w_2x^2 + w_3x^3 + \epsilon$$



Linear



Cubic

Plots show $t_{\text{new}} \pm \text{var}\{t_{\text{new}}\}$. (Black line is truth).



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

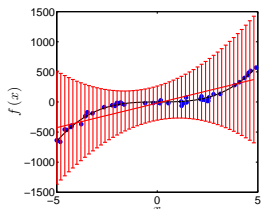
Likelihood for model selection

Summary

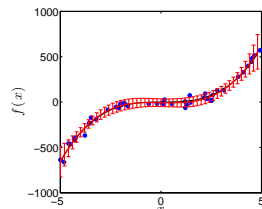
Example

Data sampled from a 3rd order polynomial function:

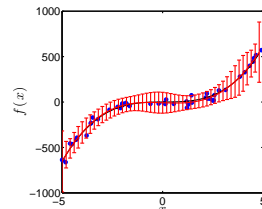
$$t = w_0 + w_1x + w_2x^2 + w_3x^3 + \epsilon$$



Linear



Cubic



6th order

Plots show $t_{\text{new}} \pm \text{var}\{t_{\text{new}}\}$. (Black line is truth).



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

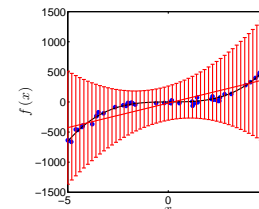
Likelihood for model selection

Summary

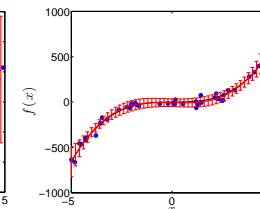
Example

Data sampled from a 3rd order polynomial function:

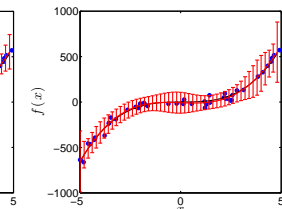
$$t = w_0 + w_1x + w_2x^2 + w_3x^3 + \epsilon$$



Linear



Cubic



6th order

Plots show $t_{\text{new}} \pm \text{var}\{t_{\text{new}}\}$. (Black line is truth).

Why does the predictive variance increase above and below the correct order?



Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

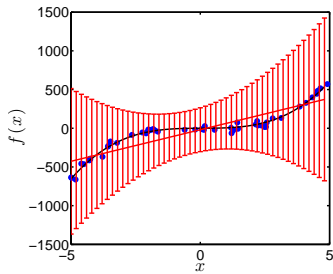
Likelihood for model selection

Summary

Not complex enough model – more ‘noise’

In practice we don't know σ^2 so substitute $\widehat{\sigma^2}$:

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

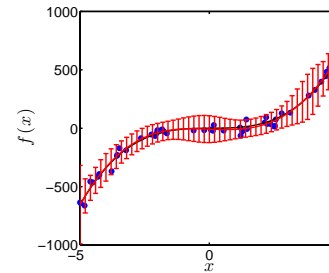


- ▶ The model is too simple.
- ▶ Some true variability can only be modelled noise.
- ▶ $\widehat{\sigma^2}$ is significantly over-estimated.
- ▶ Results in high $\text{var}\{t_{\text{new}}\}$.

Too complex model – parameters not well defined

Similarly, we substitute $\widehat{\sigma^2}$ into expression for $\text{cov}\{\widehat{\mathbf{w}}\}$:

$$\text{cov}\{\widehat{\mathbf{w}}\} = \widehat{\sigma^2} (\mathbf{X}^T \mathbf{X})^{-1}$$

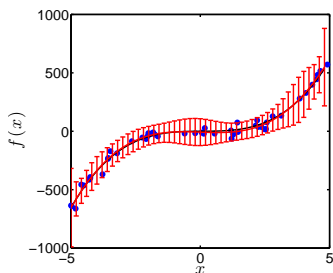


- ▶ 6th order model is too flexible.
- ▶ Many sets of parameters lead to a good model.
- ▶ Means that $\text{cov}\{\widehat{\mathbf{w}}\}$ is high.

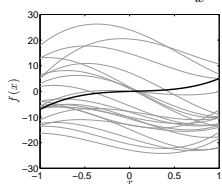
Too complex model – parameters not well defined

Similarly, we substitute $\widehat{\sigma^2}$ into expression for $\text{cov}\{\widehat{\mathbf{w}}\}$:

$$\text{cov}\{\widehat{\mathbf{w}}\} = \widehat{\sigma^2} (\mathbf{X}^T \mathbf{X})^{-1}$$



- ▶ 6th order model is too flexible.
- ▶ Many sets of parameters lead to a good model.
- ▶ Means that $\text{cov}\{\widehat{\mathbf{w}}\}$ is high.

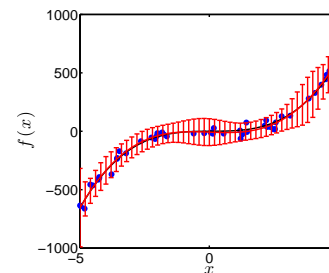


- ▶ ‘good’ 6th order models.

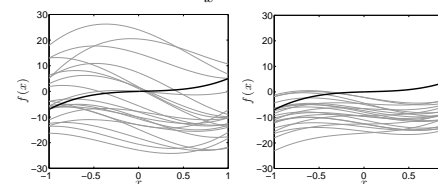
Too complex model – parameters not well defined

Similarly, we substitute $\widehat{\sigma^2}$ into expression for $\text{cov}\{\widehat{\mathbf{w}}\}$:

$$\text{cov}\{\widehat{\mathbf{w}}\} = \widehat{\sigma^2} (\mathbf{X}^T \mathbf{X})^{-1}$$



- ▶ 6th order model is too flexible.
- ▶ Many sets of parameters lead to a good model.
- ▶ Means that $\text{cov}\{\widehat{\mathbf{w}}\}$ is high.



- ▶ ‘good’ 6th order models.
- ▶ ‘good’ 3rd order models.

Olympic prediction

Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

Prediction

Likelihood for model selection

Summary

Linear model:

$$t = w_0 + w_1 X + \epsilon$$



Olympic prediction

Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

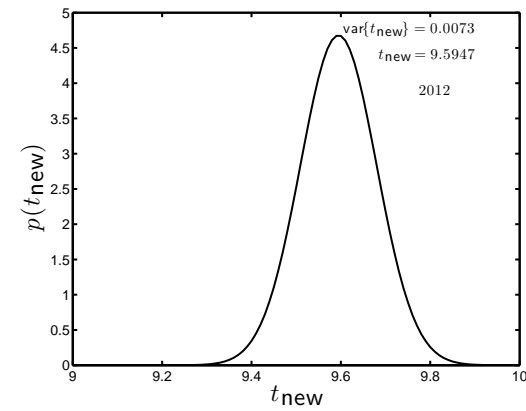
Prediction

Likelihood for model selection

Summary

Linear model:

$$t = w_0 + w_1 X + \epsilon$$



Olympic prediction

Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

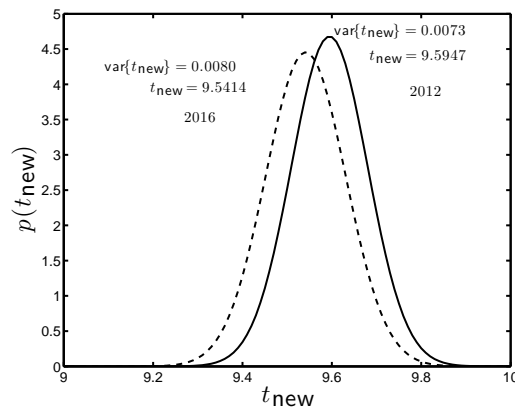
Prediction

Likelihood for model selection

Summary

Linear model:

$$t = w_0 + w_1 X + \epsilon$$



Predictive variance increases as we get further from the training data.



Olympic prediction

Introduction

D. Dubhashi

Confidence in parameter estimates

Predictions

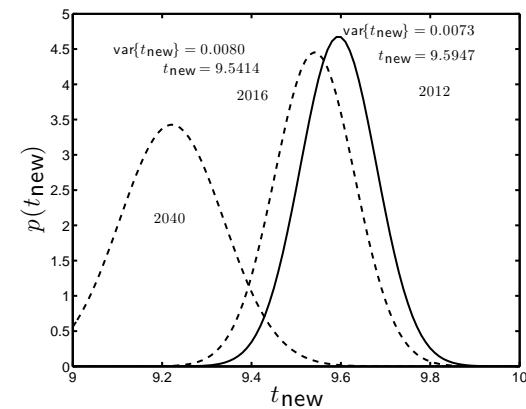
Prediction

Likelihood for model selection

Summary

Linear model:

$$t = w_0 + w_1 X + \epsilon$$



Predictive variance increases as we get further from the training data.



Summary

- ▶ Decided to model the noise.
- ▶ Recapped random variables.
- ▶ Introduced likelihood and maximised it to find $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$.
- ▶ What did it buy us?

- Introduction
- D. Dubhashi
- Confidence in parameter estimates
- Predictions
- Prediction
- Likelihood for model selection
- Summary

Summary

- ▶ Decided to model the noise.
- ▶ Recapped random variables.
- ▶ Introduced likelihood and maximised it to find $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$.
- ▶ What did it buy us?
- ▶ We can now:
 - ▶ Quantify the uncertainty in our parameters.
 - ▶ Quantify the uncertainty in our predictions.
 - ▶ This is very important in all applications....

- Introduction
- D. Dubhashi
- Confidence in parameter estimates
- Predictions
- Prediction
- Likelihood for model selection
- Summary

Summary

- ▶ Decided to model the noise.
- ▶ Recapped random variables.
- ▶ Introduced likelihood and maximised it to find $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$.
- ▶ What did it buy us?
- ▶ We can now:
 - ▶ Quantify the uncertainty in our parameters.
 - ▶ Quantify the uncertainty in our predictions.
 - ▶ This is very important in all applications....
- ▶ What next?
 - ▶ Going Bayesian.
 - ▶ Got to forget about single parameter values - parameters are random variables too.

- Introduction
- D. Dubhashi
- Confidence in parameter estimates
- Predictions
- Prediction
- Likelihood for model selection
- Summary

Aside - from one model to many

- ▶ All of our efforts so far have been to find the 'best' model:
 - ▶ The one that minimises the loss.
 - ▶ The one that maximises the likelihood.
- ▶ Given the uncertainty, maybe we shouldn't trust one on its own?
- ▶ Consider the following RV:

$$p(\mathbf{q}) = \mathcal{N}(\hat{\mathbf{w}}, \text{cov}\{\hat{\mathbf{w}}\})$$

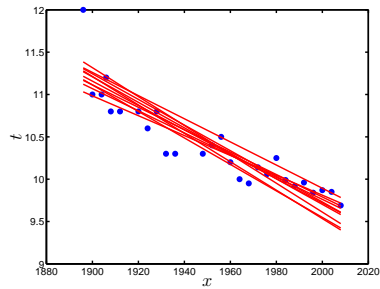
- ▶ Samples of this RV \mathbf{q}_s are **models** (assume $\hat{\sigma}^2$ is fixed)
- ▶ We can generate lots of good models...

- Introduction
- D. Dubhashi
- Confidence in parameter estimates
- Predictions
- Prediction
- Likelihood for model selection
- Summary

- ▶ Sample lots of \mathbf{q} from:

$$p(\mathbf{q}) = \mathcal{N}(\hat{\mathbf{w}}, \text{cov}\{\hat{\mathbf{w}}\})$$

- ▶ Each corresponds to a model.

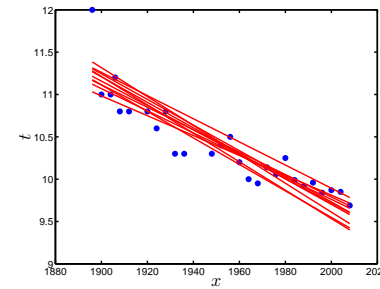


- ▶ Sample lots of \mathbf{q} from:

$$p(\mathbf{q}) = \mathcal{N}(\hat{\mathbf{w}}, \text{cov}\{\hat{\mathbf{w}}\})$$

- ▶ Each corresponds to a model.
- ▶ Compute a prediction from each one:

$$t_s = \mathbf{q}_s^T \mathbf{x}_{\text{new}}$$



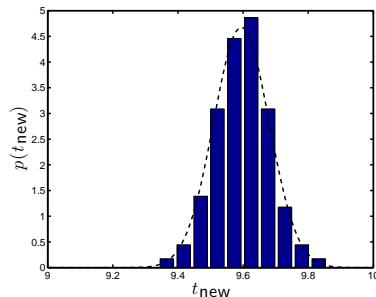
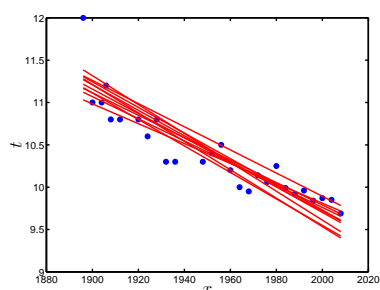
- ▶ Sample lots of \mathbf{q} from:

$$p(\mathbf{q}) = \mathcal{N}(\hat{\mathbf{w}}, \text{cov}\{\hat{\mathbf{w}}\})$$

- ▶ Each corresponds to a model.
- ▶ Compute a prediction from each one:

$$t_s = \mathbf{q}_s^T \mathbf{x}_{\text{new}}$$

- ▶ Look at the distribution of predictions:



Do we need to take samples at all?

- ▶ Take an expectation...

$$\mathbf{E}_{p(\mathbf{q})} \{t_{\text{new}}\} = \int t_{\text{new}} \mathcal{N}(\hat{\mathbf{w}}, \text{cov}\{\hat{\mathbf{w}}\}) dt_{\text{new}}$$

- ▶ We'll see more of this in the next lecture....

