

Probabilistic Approach to Linear Regression

Morteza H. Chehreghani

`morteza.chehreghani@chalmers.se`

Department of Computer Science and Engineering
Chalmers University of Technology

January 23, 2024

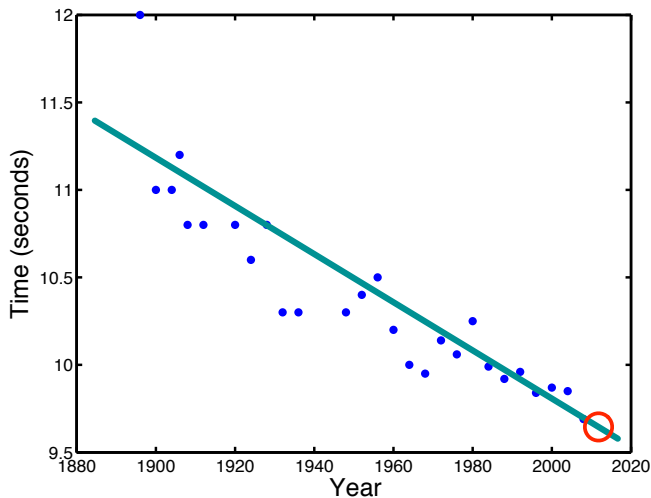
Reference

The content and the slides are adapted from

S. Rogers and M. Girolami, A First Course in Machine Learning (FCML), 2nd edition, Chapman & Hall/CRC 2016, ISBN: 9781498738484

Some data and a problem

Use the model (line) to *predict* the winning time in 2012.



Recipe for a linear model

$$t_n = w_0 + w_1x_{n,1} + w_2x_{n,2} + w_3x_{n,3} + \dots + w_Dx_{n,D}$$


$$\mathbf{x}_n = \begin{bmatrix} 1 \\ x_{n,1} \\ x_{n,2} \\ \vdots \\ x_{n,D} \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,D} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & x_{N,2} & \dots & x_{N,D} \end{bmatrix} \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix},$$

Recipe for a linear model

$$t_n = w_0 + w_1 x_{n,1} + w_2 x_{n,2} + w_3 x_{n,3} + \dots + w_D x_{n,D}$$

$$\mathbf{x}_n = \begin{bmatrix} 1 \\ x_{n,1} \\ x_{n,2} \\ \vdots \\ x_{n,D} \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,D} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & x_{N,2} & \dots & x_{N,D} \end{bmatrix} \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix},$$

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}, \quad \text{Model : } t_n = \mathbf{w}^T \mathbf{x}_n, \quad \text{or} \quad \mathbf{t} = \mathbf{X} \mathbf{w}$$



What about the errors?

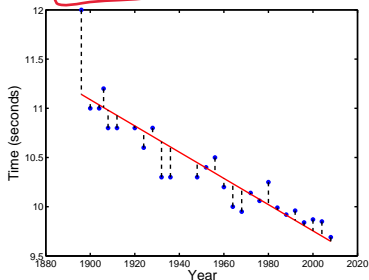
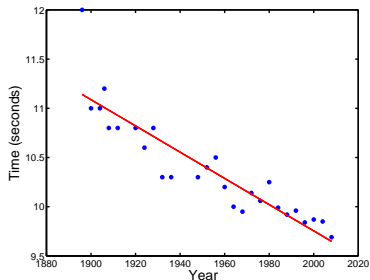
$$t_n = w_0 + w_1 x_{n,1} + w_2 x_{n,2} + w_3 x_{n,3} + \dots + w_D x_{n,D} = \sum_{d=0}^D w_d x_{n,d} = \mathbf{w}^T \mathbf{x}_n$$

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N \left(t_n - \mathbf{w}^T \mathbf{x}_n \right)^2 = \frac{1}{N} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})$$

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

$\left. \begin{matrix} w_0 \\ w_1 \end{matrix} \right\}$

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

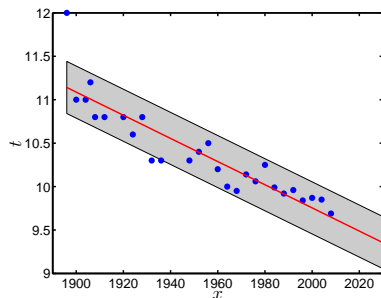


We **should** model the errors

- ▶ We know they're there - shouldn't ignore them.

We **should** model the errors

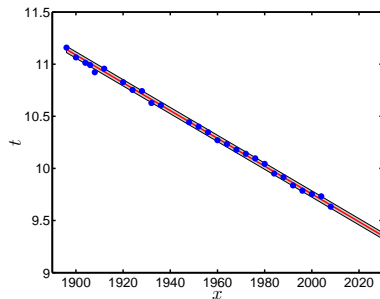
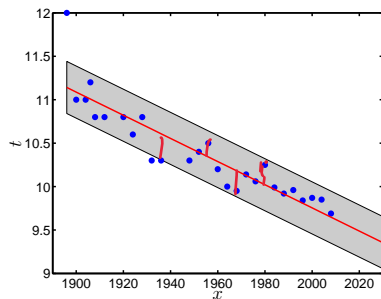
- ▶ We know they're there - shouldn't ignore them.
- ▶ They tell us how confident our predictions should be:



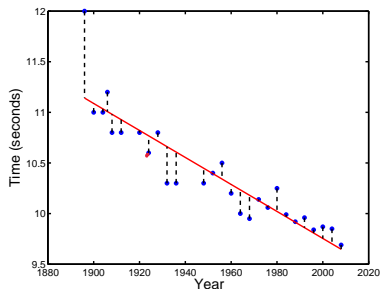
We **should** model the errors

- ▶ We know they're there - shouldn't ignore them.

- ▶ They tell us how confident our predictions should be:



Additive errors

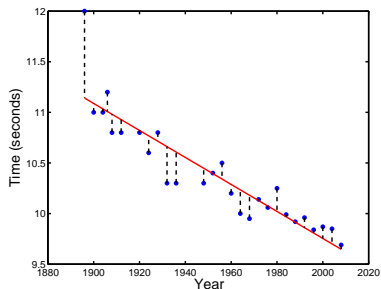


$$t_n = \mathbf{w}^T \mathbf{x}_n$$

We'll assume that the noise is an additive term in the model:

$$t_n = \underbrace{\mathbf{w}^T \mathbf{x}_n}_{w_0 + w_1 x_n} + \underbrace{\epsilon_n}_{\text{noise}}$$

Additive errors

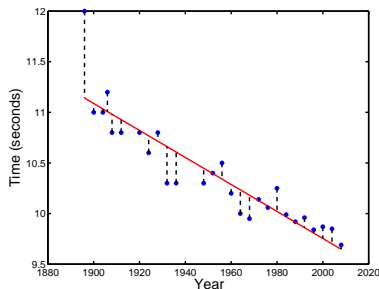


We'll assume that the noise is an additive term in the model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

What assumptions can we make about ϵ_n ?

Additive errors



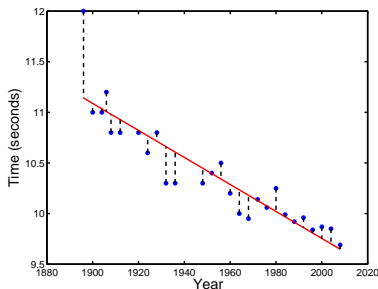
We'll assume that the noise is an additive term in the model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

What assumptions can we make about ϵ_n ?

- It's different for each n .

Additive errors



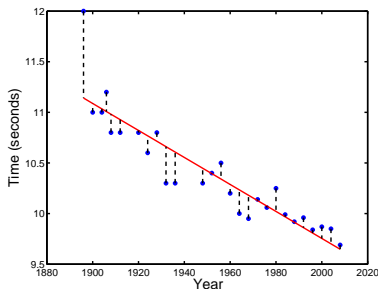
We'll assume that the noise is an additive term in the model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

What assumptions can we make about ϵ_n ?

- ▶ It's different for each n .
- ▶ It's positive and negative.

Additive errors



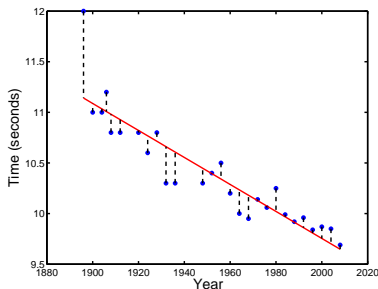
We'll assume that the noise is an additive term in the model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

What assumptions can we make about ϵ_n ?

- ▶ It's different for each n .
- ▶ It's positive and negative.
- ▶ There doesn't seem to be any relationship between ϵ at different n .

Additive errors



We'll assume that the noise is an additive term in the model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

What assumptions can we make about ϵ_n ?

- ▶ It's different for each n .
- ▶ It's positive and negative. $\rightarrow E[\epsilon_n] = 0$
- ▶ There doesn't seem to be any relationship between ϵ at different n . \rightarrow independence of t_n 's
- ▶ Looks very hard to model exactly (if it were, it wouldn't be noise!)

Gaussian noise model

- Our model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

Gaussian noise model

- ▶ Our model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

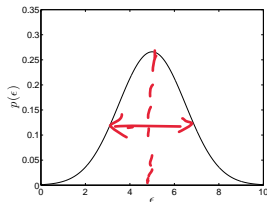
- ▶ ϵ_n is continuous.
- ▶ We need to choose $p(\epsilon)$.

Gaussian noise model

- ▶ Our model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

- ▶ ϵ_n is continuous.
- ▶ We need to choose $p(\epsilon)$.
- ▶ Gaussian:



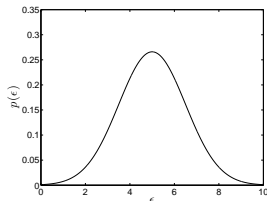
$$p(\epsilon|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2} (\epsilon - \mu)^2 \right\}$$

Gaussian noise model

- ▶ Our model:

$$t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

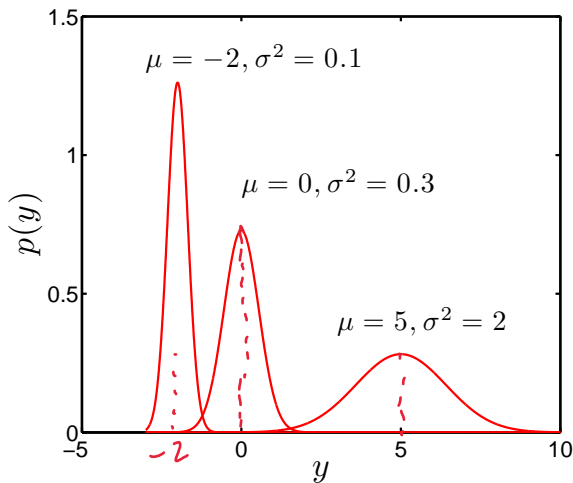
- ▶ ϵ_n is continuous.
- ▶ We need to choose $p(\epsilon)$.
- ▶ Gaussian:



$$p(\epsilon|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2} (\epsilon - \mu)^2 \right\}$$

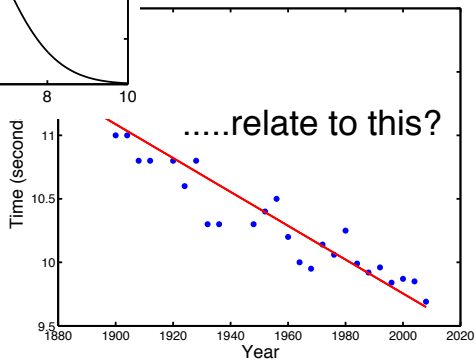
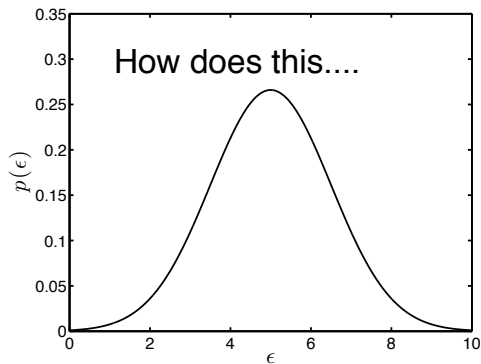
- ▶ 2 parameters: Mean μ and Variance σ^2 .

Gaussian examples



Effect of varying the mean (μ) and variance (σ^2) parameters of the Gaussian.

Generating data

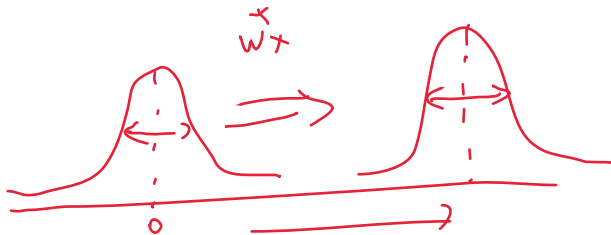


Likelihood

$$t = \overbrace{w^T x}^{\mu = w^T x, \sigma^2} + \underbrace{\epsilon}_Z \rightarrow \mu = 0, \sigma^2$$

► t is a random variable too!

$$p(t|\mathbf{x}_n, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2)$$



Likelihood

- ▶ t is a random variable too!

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

- ▶ Evaluate the density at $t = t_n$!
- ▶ At $t = t_n$ it is called the **Likelihood**, i.e., the quantity obtained when evaluating the density.
- ▶ The higher the value, the more likely t_n is given the model....

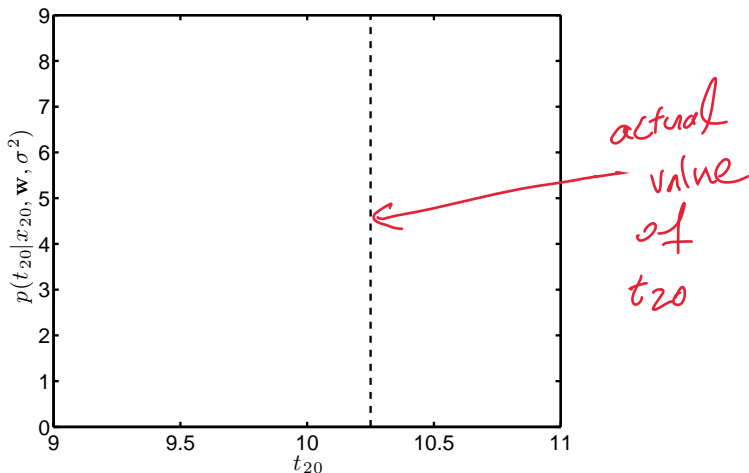
Likelihood

- ▶ t is a random variable too!

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

- ▶ Evaluate the density at $t = t_n$!
- ▶ At $t = t_n$ it is called the **Likelihood**, i.e., the quantity obtained when evaluating the density.
- ▶ The higher the value, the more likely t_n is given the model....
 - ▶the better the model is.

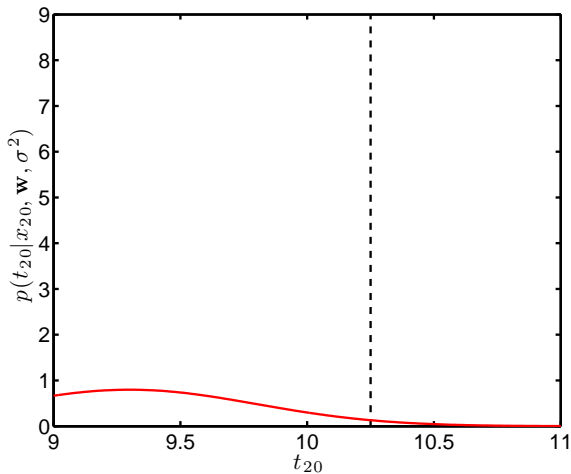
Likelihood



Lets look at the 1980 Olympics ($n = 20$).

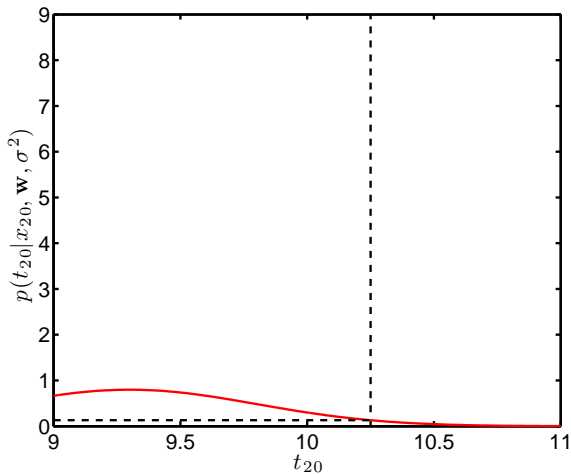
Dashed line shows t_{20} .

Likelihood



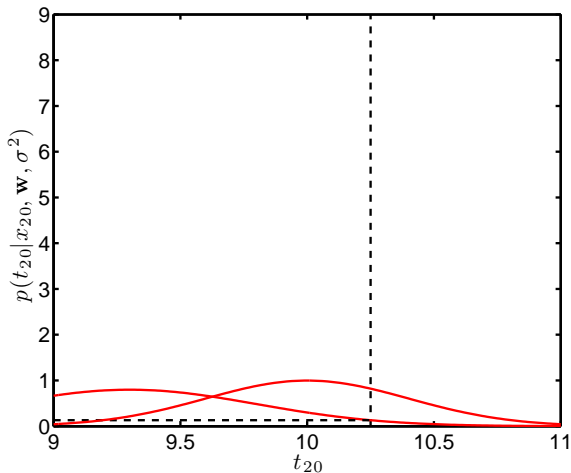
Model 1. Red line shows $\mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2)$

Likelihood



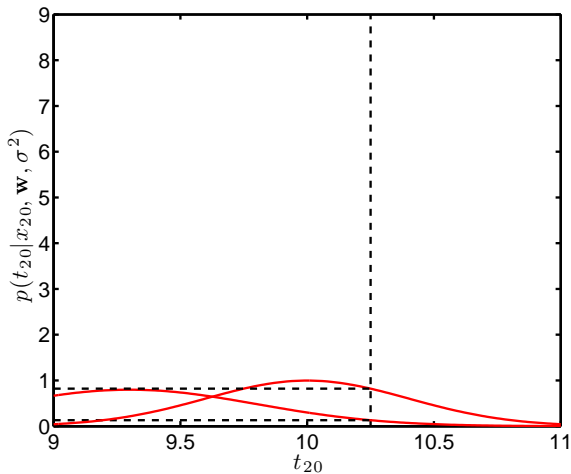
$$p(t_{20}|\dots) \approx 0.1.$$

Likelihood



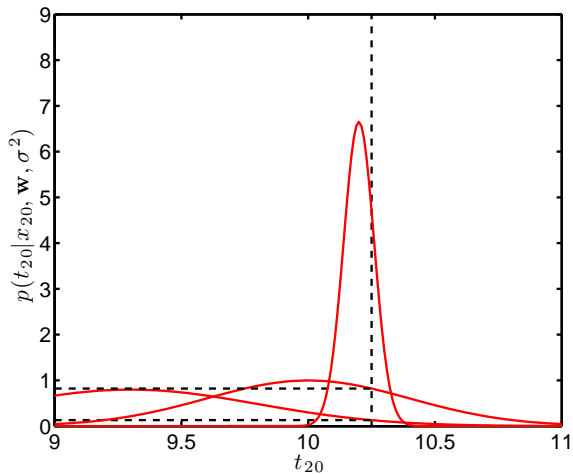
Model 2. Red line shows $\mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$ for a different \mathbf{w}

Likelihood



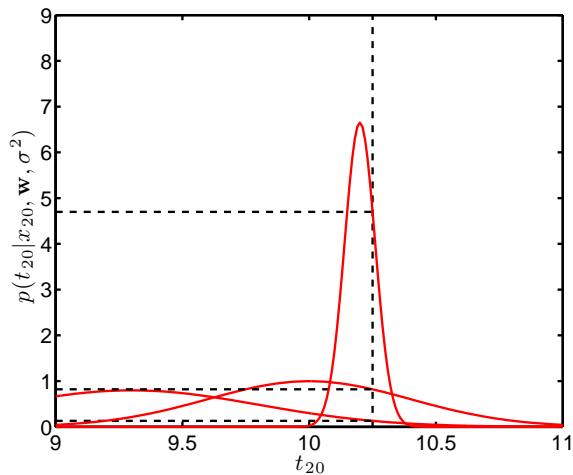
$$p(t_{20}|\dots) \approx 0.9.$$

Likelihood



Model 3.

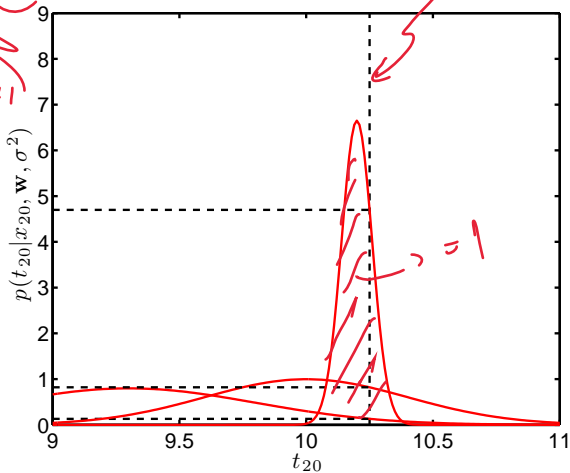
Likelihood



Model 3.

Likelihood

$$= \mathcal{N}(w, \sigma^2)$$



Model 3.

Model 3 looks best.

Likelihood

- ▶ The value we get when we evaluate the density function is called the **likelihood**.

Likelihood

- ▶ The value we get when we evaluate the density function is called the **likelihood**.
- ▶ i.e.
 - ▶ The likelihood for model 1 was 0.1.
 - ▶ The likelihood for model 2 was 0.9.
 - ▶ The likelihood for model 3 was 4.8.
- ▶ For continuous random variables, it is **not** a probability!

Likelihood

- ▶ The value we get when we evaluate the density function is called the **likelihood**.
- ▶ i.e.
 - ▶ The likelihood for model 1 was 0.1.
 - ▶ The likelihood for model 2 was 0.9.
 - ▶ The likelihood for model 3 was 4.8.
- ▶ For continuous random variables, it is **not** a probability!
- ▶ As t_n is fixed, we can find the values of \mathbf{w} and σ^2 that maximise the likelihood.
 - ▶ ...just like we found them that minimised the loss.


Likelihood optimisation

- For each input-response pair, we have a Gaussian likelihood:

$$p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2) = \mathcal{N}(\underbrace{\mathbf{w}^\top \mathbf{x}_n}_{\text{mean}}, \underbrace{\sigma^2}_{\text{variance}})$$

Likelihood optimisation

- ▶ For each input-response pair, we have a Gaussian likelihood:


$$p(t_n|\mathbf{w}, \mathbf{x}_n, \sigma^2) = \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2)$$

- ▶ To combine them all, we want the joint likelihood:

entire
dataset



$$p(t_1, \dots, t_N | \mathbf{w}, \sigma^2, \mathbf{x}_1, \dots, \mathbf{x}_N)$$

Likelihood optimisation

- ▶ For each input-response pair, we have a Gaussian likelihood:

$$p(t_n|\mathbf{w}, \mathbf{x}_n, \sigma^2) = \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2)$$

- ▶ To combine them all, we want the joint likelihood:

$$p(t_1, \dots, t_N | \mathbf{w}, \sigma^2, \mathbf{x}_1, \dots, \mathbf{x}_N)$$

- ▶ Assume that the t_n 's are independent:

$$p(t_1, \dots, t_N | \mathbf{w}, \sigma^2, \mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{n=1}^N p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2)$$

Likelihood optimisation

Finding the parameters that maximise the likelihood is expressed mathematically as:

$$\operatorname{argmax}_{\mathbf{w}, \sigma^2} \prod_{n=1}^N p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2)$$

Likelihood optimisation



Finding the parameters that maximise the likelihood is expressed mathematically as:

$$\operatorname{argmax}_{\mathbf{w}, \sigma^2} \prod_{n=1}^N p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2)$$

In fact, we'll optimise the (natural) log likelihood because it's easier.

- ▶ If we increase z , $\log(z)$ increases, if we decrease z , $\log(z)$ decreases. So, at a maximum of z , $\log(z)$ will also be at a maximum.

Likelihood optimisation

Finding the parameters that maximise the likelihood is expressed mathematically as:

$$\operatorname{argmax}_{\mathbf{w}, \sigma^2} \prod_{n=1}^N p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2)$$

In fact, we'll optimise the (natural) log likelihood because it's easier.

- If we increase z , $\log(z)$ increases, if we decrease z , $\log(z)$ decreases. So, at a maximum of z , $\log(z)$ will also be at a maximum.

$$\begin{aligned} \operatorname{argmax}_{\mathbf{w}, \sigma^2} \log \prod_{n=1}^N p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2) \\ = \sum_{n=1}^N \log p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2) \end{aligned}$$

Some re-arranging...

$$p(t_n|\mathbf{w}, \mathbf{x}_n, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(t_n - \mathbf{w}^\top \mathbf{x}_n)^2\right\}$$

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

Some re-arranging...

$$\begin{aligned} p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2) &= \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2} (t_n - \mathbf{w}^\top \mathbf{x}_n)^2 \right\} \\ \log L &= \log \prod_{n=1}^N p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2) \\ &= \sum_{n=1}^N \log p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2) \\ &= \sum_{n=1}^N \log \left(\frac{1}{\sigma\sqrt{2\pi}} \right) - \sum_{n=1}^N \frac{1}{2\sigma^2} (t_n - \mathbf{w}^\top \mathbf{x}_n)^2 \\ &= -N \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2 \end{aligned}$$

Looks familiar!

Some re-arranging...

$$\begin{aligned}p(t_n|\mathbf{w}, \mathbf{x}_n, \sigma^2) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(t_n - \mathbf{w}^\top \mathbf{x}_n)^2\right\} \\ \log L &= \log \prod_{n=1}^N p(t_n|\mathbf{w}, \mathbf{x}_n, \sigma^2) \\ &= \sum_{n=1}^N \log p(t_n|\mathbf{w}, \mathbf{x}_n, \sigma^2) \\ &= \sum_{n=1}^N \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \sum_{n=1}^N \frac{1}{2\sigma^2}(t_n - \mathbf{w}^\top \mathbf{x}_n)^2 \\ &= -N \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2\end{aligned}$$

Looks familiar! To continue (good exercise):

$$\frac{\partial \log L}{\partial \mathbf{w}} = 0, \quad \frac{\partial \log L}{\partial \sigma^2} = 0$$

A shortcut

The multi-variate Gaussian

random vector $\leftarrow \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{K/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\}$$

$K(= 2)$ is number of variables, $|\boldsymbol{\Sigma}|$ is the determinant.

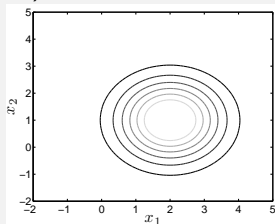
A shortcut

The multi-variate Gaussian

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, \quad p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{K/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\}$$

$K(= 2)$ is number of variables, $|\boldsymbol{\Sigma}|$ is the determinant.



$$\boldsymbol{\mu} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

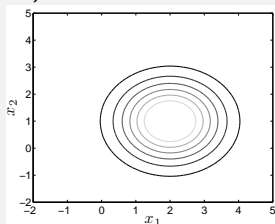
A shortcut

The multi-variate Gaussian

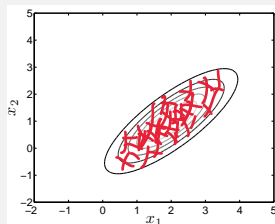
$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, \quad p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{K/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\}$$

$K(= 2)$ is number of variables, $|\boldsymbol{\Sigma}|$ is the determinant.



$$\boldsymbol{\mu} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

A shortcut

The multi-variate Gaussian

A special case:

$$\prod_{n=1}^N \mathcal{N}(\mu_n, \sigma^2) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

independence of r.v.'s

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_N \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \sigma^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}$$

*[1 0 0]
0 1 0
0 0 1]*

So, in our model:

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

Maximising the multi-variate log-likelihood

- Partial derivative w.r.t. \mathbf{w} , set to zero and solve:

$$\begin{aligned}\log L &= \log \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}) \\ \frac{\partial \log L}{\partial \mathbf{w}} &= -\frac{1}{2\sigma^2} (2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{t}) = 0 \\ \mathbf{w} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}\end{aligned}$$

Handwritten notes: A red circle around $\mathbf{X}\mathbf{w}$ in the first equation, with a red arrow pointing to \mathbf{w} and \mathbf{X} written above it. Red slashes are present under the 1 and 2 in the second equation.

Maximising the multi-variate log-likelihood

- ▶ Partial derivative w.r.t. \mathbf{w} , set to zero and solve:

$$\begin{aligned}\log L &= \log \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}) \\ \frac{\partial \log L}{\partial \mathbf{w}} &= -\frac{1}{2\sigma^2} (2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{t}) = 0 \\ \mathbf{w} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}\end{aligned}$$

- ▶ This is the same expression we've seen before!

Maximising the multi-variate log-likelihood

- ▶ Partial derivative w.r.t. \mathbf{w} , set to zero and solve:

$$\begin{aligned}\log L &= \log \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}) \\ \frac{\partial \log L}{\partial \mathbf{w}} &= -\frac{1}{2\sigma^2} (2\mathbf{X}^\top \mathbf{X}\mathbf{w} - 2\mathbf{X}^\top \mathbf{t}) = 0 \\ \mathbf{w} &= \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}}\end{aligned}$$

- ▶ This is the same expression we've seen before!
- ▶ Same for σ^2 :

$$\begin{aligned}\frac{\partial \log L}{\partial \sigma^2} &= -\frac{N}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) = 0 \\ \sigma^2 &= \underbrace{\frac{1}{N} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w})}\end{aligned}$$

Optimum parameters

- Compute optimum $\hat{\mathbf{w}}$ from:

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

- Use this to compute optimum $\hat{\sigma}^2$ from:

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

Optimum parameters

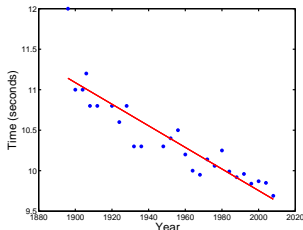
- Compute optimum $\hat{\mathbf{w}}$ from:

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

- Use this to compute optimum $\hat{\sigma}^2$ from:

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

- e.g. Olympic 100 m data (again!)



$$\hat{\mathbf{w}} = \begin{bmatrix} 36.416 \\ -0.0133 \end{bmatrix}, \quad \hat{\sigma}^2 = 0.0503$$

Confidence in parameter estimates

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

Confidence in parameter estimates

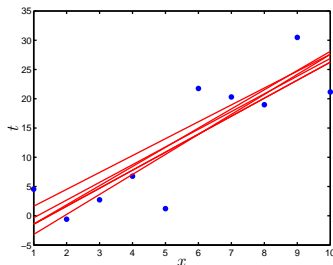
- ▶ Imagine there are **true** parameters, \mathbf{w} and σ^2 .
- ▶ How good are our estimates $\hat{\mathbf{w}}$ and $\hat{\sigma}^2$?
 - ▶ Are they correct (on average)?
 - ▶ If we could keep adding data, would we converge on the true value?

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?
- ▶ How confident should we be in our estimates?
 - ▶ Could we change parameters a little bit and still have a good model?

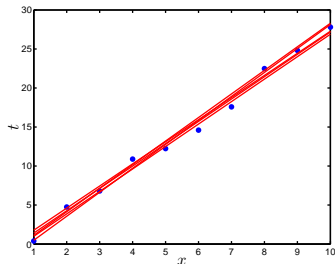
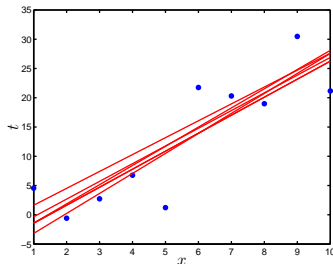
Confidence in parameter estimates

- ▶ 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?
- ▶ How confident should we be in our estimates?
 - ▶ Could we change parameters a little bit and still have a good model?



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?
- ▶ How confident should we be in our estimates?
 - ▶ Could we change parameters a little bit and still have a good model?

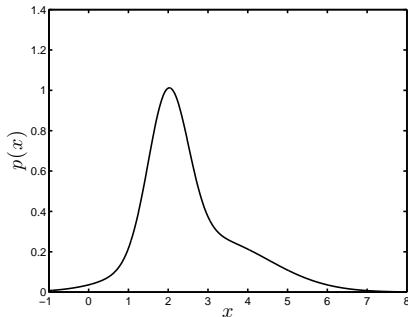


Expectations – refresher

- ▶ To progress we need to understand **Expectations**

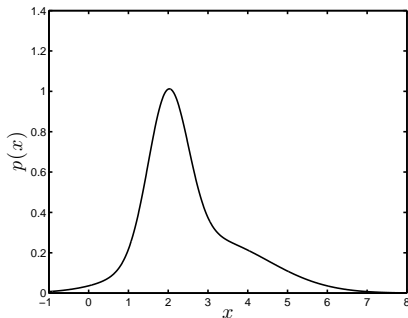
Expectations – refresher

- ▶ To progress we need to understand **Expectations**
- ▶ Imagine a random variable X with density $p(x)$



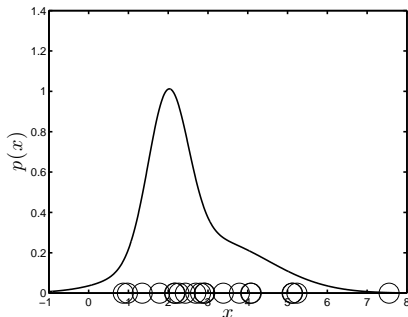
Expectations – refresher

- ▶ To progress we need to understand **Expectations**
- ▶ Imagine a random variable X with density $p(x)$
- ▶ We want to work out the average value of X , \tilde{x} .



Expectations – refresher

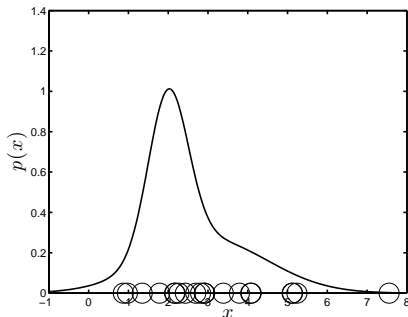
- ▶ To progress we need to understand **Expectations**
- ▶ Imagine a random variable X with density $p(x)$
- ▶ We want to work out the average value of X , \tilde{x} .
- ▶ Generate S samples, x_1, \dots, x_S



Expectations – refresher

- ▶ To progress we need to understand **Expectations**
- ▶ Imagine a random variable X with density $p(x)$
- ▶ We want to work out the average value of X , \tilde{x} .
- ▶ Generate S samples, x_1, \dots, x_S
- ▶ Average the samples:

$$\tilde{x} \approx \frac{1}{S} \sum_{s=1}^S x_s$$



Expectations – refresher

- ▶ Our sample based approximation to \tilde{x} will get better as we take more samples.

Expectations – refresher

- ▶ Our sample based approximation to \tilde{x} will get better as we take more samples.
- ▶ We can also (sometimes) compute it exactly using **expectations**.

Expectations – refresher

- ▶ Our sample based approximation to \tilde{x} will get better as we take more samples.
- ▶ We can also (sometimes) compute it exactly using **expectations**.
 - ▶ Discrete: $\tilde{x} = \mathbf{E}_{p(x)} \{x\} = \sum_x x p(x)$
- ▶ Example:
 - ▶ X is outcome of rolling die. $P(X = x) = 1/6$
 - ▶ $\tilde{x} = \sum_x x P(X = x) = 3.5$

↪ 1, 2, 3, 4, 5, 6

Expectations – refresher

- ▶ Our sample based approximation to \tilde{x} will get better as we take more samples.
- ▶ We can also (sometimes) compute it exactly using **expectations**.
 - ▶ Discrete: $\tilde{x} = \mathbf{E}_{p(x)} \{x\} = \sum_x xp(x)$
 - ▶ Continuous: $\tilde{x} = \mathbf{E}_{p(x)} \{x\} = \int_x xp(x) dx$
- ▶ Example:
 - ▶ X is outcome of rolling die. $P(X = x) = 1/6$
 - ▶ $\tilde{x} = \sum_x xP(X = x) = 3.5$
- ▶ Example:
 - ▶ X is uniform distributed RV between a and b
 - ▶ $\tilde{x} = \int_{x=a}^{x=b} xp(x) dx = (b + a)/2$

$\rightarrow p(x) = \frac{1}{b-a}$

Expectations – refresher

- ▶ In general:

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

Expectations – refresher

- ▶ In general:

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

- ▶ Some important things:

- ▶ $\mathbf{E}_{p(x)} \{f(x)\} \neq f(\mathbf{E}_{p(x)} \{x\})$
- ▶ $\mathbf{E}_{p(x)} \{kf(x)\} = k\mathbf{E}_{p(x)} \{f(x)\}$

Expectations – refresher

- ▶ In general:

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

- ▶ Some important things:

- ▶ $\mathbf{E}_{p(x)} \{f(x)\} \neq f(\mathbf{E}_{p(x)} \{x\})$
- ▶ $\mathbf{E}_{p(x)} \{kf(x)\} = k\mathbf{E}_{p(x)} \{f(x)\}$
- ▶ Mean, $\mu = \mathbf{E}_{p(x)} \{x\}$

Expectations – refresher

- ▶ In general:

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

- ▶ Some important things:

- ▶ $\mathbf{E}_{p(x)} \{f(x)\} \neq f(\mathbf{E}_{p(x)} \{x\})$
- ▶ $\mathbf{E}_{p(x)} \{kf(x)\} = k\mathbf{E}_{p(x)} \{f(x)\}$
- ▶ Mean, $\mu = \mathbf{E}_{p(x)} \{x\}$
- ▶ Variance: $\sigma^2 = \mathbf{E}_{p(x)} \{(x - \mu)^2\} = \underbrace{\mathbf{E}_{p(x)} \{x^2\}}_{\neq \mu^2} - \underbrace{(\mathbf{E}_{p(x)} \{x\})^2}_{\mu^2}$

Expectations – refresher

- ▶ In general:

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

- ▶ Some important things:

- ▶ $\mathbf{E}_{p(x)} \{f(x)\} \neq f(\mathbf{E}_{p(x)} \{x\})$

- ▶ $\mathbf{E}_{p(x)} \{kf(x)\} = k\mathbf{E}_{p(x)} \{f(x)\}$

- ▶ Mean, $\mu = \mathbf{E}_{p(x)} \{x\}$

- ▶ Variance: $\sigma^2 = \mathbf{E}_{p(x)} \{(x - \mu)^2\} = \mathbf{E}_{p(x)} \{x^2\} - (\mathbf{E}_{p(x)} \{x\})^2$

- ▶ For vectors of random variables:

- ▶ $\mathbf{E}_{p(\mathbf{x})} \{f(\mathbf{x})\} = \int f(\mathbf{x})p(\mathbf{x}) d\mathbf{x}$

Expectations – refresher

- ▶ In general:

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

- ▶ Some important things:

- ▶ $\mathbf{E}_{p(x)} \{f(x)\} \neq f(\mathbf{E}_{p(x)} \{x\})$
- ▶ $\mathbf{E}_{p(x)} \{kf(x)\} = k\mathbf{E}_{p(x)} \{f(x)\}$
- ▶ Mean, $\mu = \mathbf{E}_{p(x)} \{x\}$
- ▶ Variance: $\sigma^2 = \mathbf{E}_{p(x)} \{(x - \mu)^2\} = \mathbf{E}_{p(x)} \{x^2\} - (\mathbf{E}_{p(x)} \{x\})^2$

- ▶ For vectors of random variables:

- ▶ $\mathbf{E}_{p(\mathbf{x})} \{f(\mathbf{x})\} = \int f(\mathbf{x})p(\mathbf{x}) d\mathbf{x}$
- ▶ Mean: $\boldsymbol{\mu} = \mathbf{E}_{p(\mathbf{x})} \{\mathbf{x}\}$
- ▶ Covariance:

$$\begin{aligned} \text{cov}\{\mathbf{x}\} &= \mathbf{E}_{p(\mathbf{x})} \{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T\} \\ &= \mathbf{E}_{p(\mathbf{x})} \{\mathbf{x}\mathbf{x}^T\} - \mathbf{E}_{p(\mathbf{x})} \{\mathbf{x}\} \mathbf{E}_{p(\mathbf{x})} \{\mathbf{x}^T\} \end{aligned}$$

Expectations – Gaussians

- ▶ Uni-variate

- ▶ $p(x|\mu, \sigma^2) = \mathcal{N}(\mu, \sigma^2)$

- ▶ Mean: $\mathbf{E}_{p(x)} \{x\} = \mu$

- ▶ Variance: $\mathbf{E}_{p(x)} \{(x - \mu)^2\} = \sigma^2$

Expectations – Gaussians

► Uni-variate

- $p(x|\mu, \sigma^2) = \mathcal{N}(\mu, \sigma^2)$
- Mean: $\mathbf{E}_{p(x)} \{x\} = \mu$
- Variance: $\mathbf{E}_{p(x)} \{(x - \mu)^2\} = \sigma^2$

► Multi-variate

- $p(\mathbf{x}|\mu, \sigma^2) = \mathcal{N}(\mu, \Sigma)$
- Mean: $\mathbf{E}_{p(\mathbf{x})} \{\mathbf{x}\} = \mu$
- Covariance: $\mathbf{E}_{p(\mathbf{x})} \{(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T\} = \Sigma$

Back to the model...

- ▶ Parameter estimates:

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

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

Back to the model...

- ▶ Parameter estimates:

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \\ \hat{\sigma}^2 &= \frac{1}{N} (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})\end{aligned}$$

- ▶ **True** values: \mathbf{w} , σ^2

Back to the model...

- ▶ Parameter estimates:

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \\ \hat{\sigma}^2 &= \frac{1}{N} (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})\end{aligned}$$

- ▶ **True** values: \mathbf{w} , σ^2
- ▶ Our model:

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{X} \mathbf{w}, \sigma^2 \mathbf{I})$$

Back to the model...

- ▶ Parameter estimates:

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \\ \hat{\sigma}^2 &= \frac{1}{N} (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})\end{aligned}$$

- ▶ **True** values: \mathbf{w} , σ^2
- ▶ Our model:

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{X} \mathbf{w}, \sigma^2 \mathbf{I})$$

- ▶ What's $\mathbf{E}_{p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\}$?

Back to the model...

- ▶ Parameter estimates:

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \\ \hat{\sigma}^2 &= \frac{1}{N} (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{t} - \mathbf{X} \hat{\mathbf{w}})\end{aligned}$$

- ▶ **True** values: \mathbf{w} , σ^2

- ▶ Our model:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$$

- ▶ What's $\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\}$?

- ▶ What do we expect our parameter estimate to be?

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

We'll try and find $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \hat{\mathbf{w}} \}$ in terms of the true value \mathbf{w} :

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

We'll try and find $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \hat{\mathbf{w}} \}$ in terms of the true value \mathbf{w} :

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

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

We'll try and find $\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\mathbf{w}} \}$ in terms of the true value \mathbf{w} :

$$\begin{aligned} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\mathbf{w}} \} &= \int \underline{\hat{\mathbf{w}}} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t} \\ &= \int (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \mathbf{t} \} \end{aligned}$$

$\mathbb{E}_P \{ \mathbf{t} \}$

9.

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

We'll try and find $\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\mathbf{w}} \}$ in terms of the true value \mathbf{w} :

$$\begin{aligned} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \hat{\mathbf{w}} \} &= \int \hat{\mathbf{w}} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t} \\ &= \int (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \mathbf{t} \} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w} \end{aligned}$$

$$t_n = \underbrace{w^T x_n}_{\uparrow} + \epsilon_n \quad \downarrow \quad E(\epsilon_n) = 0$$

$$E(E_n) =$$

$$\begin{array}{c} \mathbf{X} \mathbf{w} \\ \downarrow \\ \left[\begin{array}{c} \mathbf{X}^T \mathbf{X} \end{array} \right] \end{array}$$

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

We'll try and find $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \hat{\mathbf{w}} \}$ in terms of the true value \mathbf{w} :

$$\begin{aligned} \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \hat{\mathbf{w}} \} &= \int \hat{\mathbf{w}} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t} \\ &= \int (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) d\mathbf{t} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \mathbf{t} \} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{I} \mathbf{w} \\ \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \hat{\mathbf{w}} \} &= \mathbf{I} \mathbf{w} = \mathbf{w} \end{aligned}$$

$$A^{-1} A = I$$

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

We'll try and find $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \hat{\mathbf{w}} \}$ in terms of the true value \mathbf{w} :

$$\begin{aligned} \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \hat{\mathbf{w}} \} &= \int \hat{\mathbf{w}} p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2) d\mathbf{t} \\ &= \int (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2) d\mathbf{t} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \mathbf{t} \} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w} \\ \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{ \hat{\mathbf{w}} \} &= \mathbf{I} \mathbf{w} = \mathbf{w} \end{aligned}$$

bias : diff of $E(\hat{w})$, w

$\hat{\mathbf{w}}$ is unbiased

On average, we expect our estimate to equal the true value!



$\hat{w} \neq w$

$\hat{z} \rightarrow \begin{cases} -1 \\ +1 \end{cases}, z=0$

$\text{cov}\{\hat{\mathbf{w}}\}$ *vector*

- What does $\text{cov}\{\hat{\mathbf{w}}\}$ tell us?

$\text{cov}\{\hat{\mathbf{w}}\}$

- ▶ What does $\text{cov}\{\hat{\mathbf{w}}\}$ tell us?
- ▶ Recall the linear model, $\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$

$\text{cov}\{\hat{\mathbf{w}}\}$

- ▶ What does $\text{cov}\{\hat{\mathbf{w}}\}$ tell us?
- ▶ Recall the linear model, $\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$

$$\text{cov}\{\hat{\mathbf{w}}\} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

$\text{cov}\{\hat{\mathbf{w}}\}$

- ▶ What does $\text{cov}\{\hat{\mathbf{w}}\}$ tell us?

- ▶ Recall the linear model, $\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$

$$\text{cov}\{\hat{\mathbf{w}}\} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

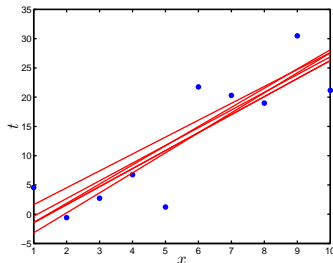
- ▶ Tells us how well defined the parameters are by the data. How much can the parameters vary and still give a **good** model.
 - ▶ a and c : how much can we change w_0 and w_1 .
 - ▶ b : how the values should be changed together.

$\text{cov}\{\hat{\mathbf{w}}\}$

- ▶ What does $\text{cov}\{\hat{\mathbf{w}}\}$ tell us?
- ▶ Recall the linear model, $\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$

$$\text{cov}\{\hat{\mathbf{w}}\} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

- ▶ Tells us how well defined the parameters are by the data. How much can the parameters vary and still give a **good** model.
 - ▶ a and c : how much can we change w_0 and w_1 .
 - ▶ b : how the values should be changed together.

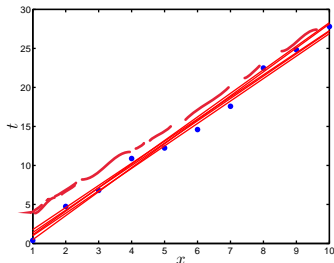
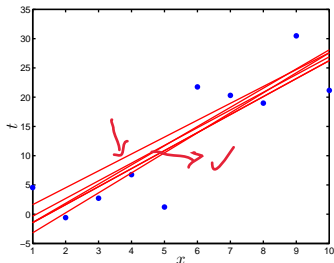


$\text{cov}\{\hat{\mathbf{w}}\}$

- ▶ What does $\text{cov}\{\hat{\mathbf{w}}\}$ tell us?
- ▶ Recall the linear model, $\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$

$$\text{cov}\{\hat{\mathbf{w}}\} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

- ▶ Tells us how well defined the parameters are by the data. How much can the parameters vary and still give a **good** model.
 - ▶ a and c : how much can we change w_0 and w_1 .
 - ▶ b : how the values should be changed together.



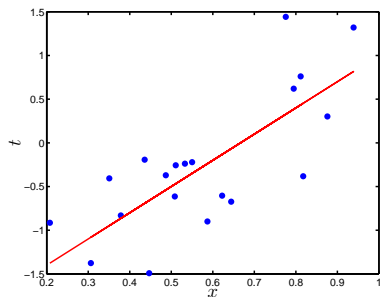
$$\text{COV}\{\hat{\mathbf{w}}\}$$

$$\begin{aligned} \text{COV}\{\hat{\mathbf{w}}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \hat{\mathbf{w}}\hat{\mathbf{w}}^T \right\} \\ &\quad - \underbrace{\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \hat{\mathbf{w}} \right\}}_{\mathbf{w}} \underbrace{\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \hat{\mathbf{w}} \right\}^T}_{\mathbf{w}^T} \end{aligned}$$

$$\text{COV}\{\hat{\mathbf{w}}\}$$

$$\begin{aligned}\text{COV}\{\hat{\mathbf{w}}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \hat{\mathbf{w}}\hat{\mathbf{w}}^T \right\} \\ &\quad - \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \hat{\mathbf{w}} \right\} \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \hat{\mathbf{w}} \right\}^T \\ &= \mathbf{E} \left\{ \hat{\mathbf{w}}\hat{\mathbf{w}}^T \right\} - \mathbf{w}\mathbf{w}^T \\ &= \vdots \\ \text{COV}\{\hat{\mathbf{w}}\} &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\end{aligned}$$

Example

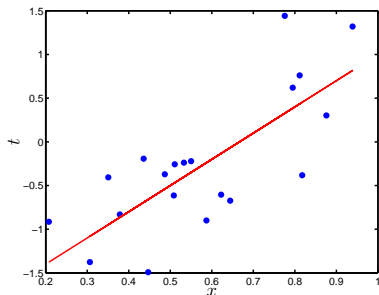


$$t_n = \underbrace{-2}_{w_0} + \underbrace{3}_{w_1}x_n + \epsilon_n$$

$$p(\epsilon_n) = \mathcal{N}(0, \sigma^2)$$

$$\sigma^2 = \underline{0.5^2}$$

Example



$$t_n = -2 + 3x_n + \epsilon_n$$

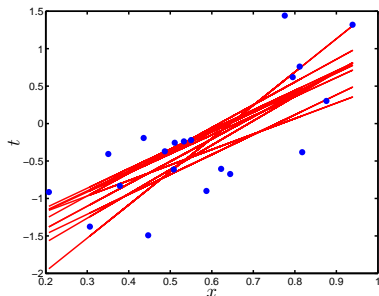
$$p(\epsilon_n) = \mathcal{N}(0, \sigma^2)$$

$$\sigma^2 = 0.5^2$$

$$\hat{\mathbf{w}} = \begin{bmatrix} -1.95 \\ 2.94 \end{bmatrix}, \text{cov}\{\hat{\mathbf{w}}\} = \begin{bmatrix} 0.1195 & -0.1847 \\ -0.1847 & 0.3190 \end{bmatrix}$$

Example

$\hat{\mathbf{w}} \rightarrow \left\{ \begin{array}{l} \text{bias} \\ \text{variance / co-var} \end{array} \right.$



$$t_n = -2 + 3x_n + \epsilon_n$$

$$p(\epsilon_n) = \mathcal{N}(0, \sigma^2)$$

$$\sigma^2 = 0.5^2$$

$$\hat{\mathbf{w}} = \begin{bmatrix} -1.95 \\ 2.94 \end{bmatrix}, \text{cov}\{\hat{\mathbf{w}}\} = \begin{bmatrix} 0.1195 & -0.1847 \\ -0.1847 & 0.3190 \end{bmatrix}$$

$\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\}$ – beyond this class

We saw that $\widehat{\mathbf{w}}$ was unbiased, what about $\widehat{\sigma^2}$?

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

Useful identity

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

Another useful identity

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

$\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\}$ – beyond this class

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

Where D is the number of columns in \mathbf{X} (the number of elements in \mathbf{w}).

Another useful identity

$$\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA})$$

$\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\}$ – beyond this class

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

- ▶ In general $D < N$.
- ▶ So $1 - D/N < 1$.
- ▶ So $\widehat{\sigma^2} < \sigma^2$

$\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\}$ – beyond this class

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

- ▶ In general $D < N$.
- ▶ So $1 - D/N < 1$.
- ▶ So $\widehat{\sigma^2} < \sigma^2$
- ▶ $\widehat{\sigma^2}$ is biased and will generally be too low.

$\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\}$ – beyond this class

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

- ▶ In general $D < N$.
- ▶ So $1 - D/N < 1$.
- ▶ So $\widehat{\sigma^2} < \sigma^2$
- ▶ $\widehat{\sigma^2}$ is biased and will generally be too low.
- ▶ Why?
 - ▶ Because it is based on $\widehat{\mathbf{w}}$ which will, in general, be closer to the data than \mathbf{w} .

$\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\} - \text{beyond this class}$

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

- ▶ In general $D < N$.
- ▶ So $1 - D/N < 1$.
- ▶ So $\widehat{\sigma^2} < \sigma^2$
- ▶ $\widehat{\sigma^2}$ is biased and will generally be too low.
- ▶ Why?
 - ▶ Because it is based on $\widehat{\mathbf{w}}$ which will, in general, be closer to the data than \mathbf{w} .
- ▶ As N increases, $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\} \rightarrow \sigma^2$

$\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\}$ – beyond this class

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

- ▶ In general $D < N$.
- ▶ So $1 - D/N < 1$.
- ▶ So $\widehat{\sigma^2} < \sigma^2$
- ▶ $\widehat{\sigma^2}$ is biased and will generally be too low.
- ▶ Why?
 - ▶ Because it is based on $\widehat{\mathbf{w}}$ which will, in general, be closer to the data than \mathbf{w} .
- ▶ As N increases, $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\} \rightarrow \sigma^2$
- ▶ To think about – what if $D = N$ or $D > N$?

Example – beyond this class

Generate 100 datasets from the following model:

$$t_n = w_0 + w_1 x_n + \epsilon_n, \quad p(\epsilon_n) = \mathcal{N}(0, 0.25)$$

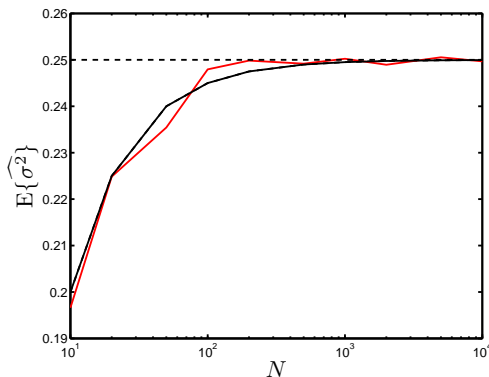
For $N = [10, 20, 50, 100, 200, 500, 1000, 2000, 5000, 10000]$

Example – beyond this class

Generate 100 datasets from the following model:

$$t_n = w_0 + w_1 x_n + \epsilon_n, \quad p(\epsilon_n) = \mathcal{N}(0, 0.25)$$

For $N = [10, 20, 50, 100, 200, 500, 1000, 2000, 5000, 10000]$



Red curve – average $\hat{\sigma}^2$ over 100 datasets. Black curve – theoretical value. Dashed line – true value.

Summary

- ▶ Computed $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{\hat{\mathbf{w}}\} = \mathbf{w}$
 - ▶ $\hat{\mathbf{w}}$ is **unbiased**.

Summary

- ▶ Computed $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)}\{\hat{\mathbf{w}}\} = \mathbf{w}$
 - ▶ $\hat{\mathbf{w}}$ is **unbiased**.
- ▶ Computed $\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$
 - ▶ Tells us how much slack there is in our parameters.

Summary

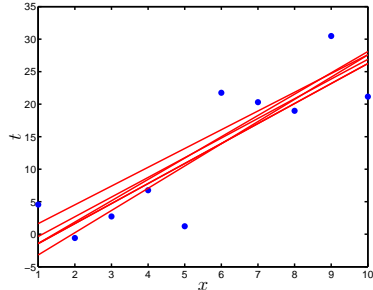
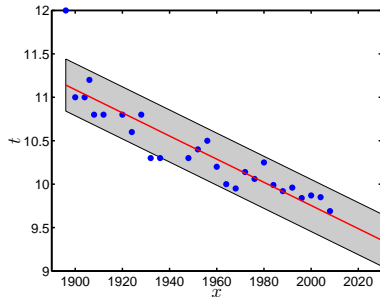
- ▶ Computed $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{\widehat{\mathbf{w}}\} = \mathbf{w}$
 - ▶ $\widehat{\mathbf{w}}$ is **unbiased**.
- ▶ Computed $\text{cov}\{\widehat{\mathbf{w}}\} = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$
 - ▶ Tells us how much slack there is in our parameters.
- ▶ Computed $\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{\widehat{\sigma^2}\} = \sigma^2(1 - D/N)$ [*beyond this class!*]
 - ▶ $\widehat{\sigma^2}$ is **biased**.
 - ▶ Gets better and better as we get more data.

Predictions

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

Predictions

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



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:

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

$$+ \cancel{\epsilon_{\text{new}}} \\ E(\epsilon_{\text{new}}) = 0$$

- ▶ Is this sensible?

< bias
variance

Predictions

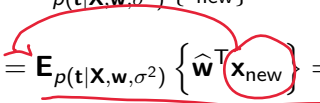
- ▶ Our model is defined as:

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

- ▶ 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!

$$E(\hat{\mathbf{w}}^T) = \mathbf{w}^T$$

Predictions

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

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

Predictions

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

$$\begin{aligned}\text{var}\{t_{\text{new}}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{t_{\text{new}}^2\} - \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \{t_{\text{new}}\}^2 \\&= \mathbf{E} \left\{ (\hat{\mathbf{w}}^T \mathbf{x}_{\text{new}})^2 \right\} - (\mathbf{w}^T \mathbf{x}_{\text{new}})^2 \\&= \mathbf{x}_{\text{new}}^T \mathbf{E} \left\{ \hat{\mathbf{w}} \hat{\mathbf{w}}^T \right\} \mathbf{x}_{\text{new}} - \mathbf{x}_{\text{new}}^T \mathbf{w} \mathbf{w}^T \mathbf{x}_{\text{new}} \\&= \vdots \\ \text{var}\{t_{\text{new}}\} &= \sigma^2 \mathbf{x}_{\text{new}}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{\text{new}}\end{aligned}$$

Prediction and variance

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

Prediction and variance

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

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

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

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:

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

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.

Example

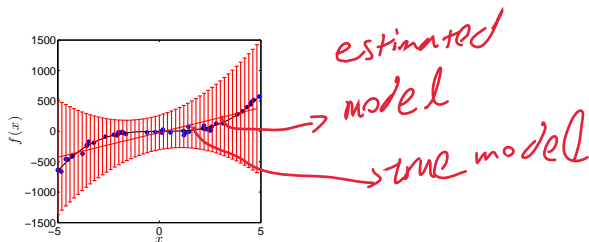
Data sampled from a 3rd order polynomial function:

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

Example

Data sampled from a 3rd order polynomial function:

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



Linear

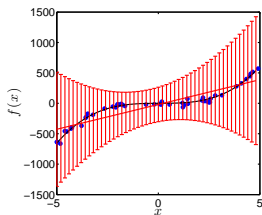
Plots show $t_{\text{new}} \pm \text{var}\{t_{\text{new}}\}$. (Black line is truth).

$w_0 + w_1x$

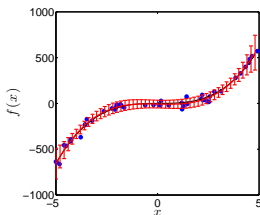
Example

Data sampled from a 3rd order polynomial function:

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



Linear



Cubic

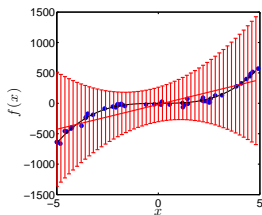
Plots show $t_{\text{new}} \pm \text{var}\{t_{\text{new}}\}$. (Black line is truth).

$w_0 + w_1x + w_2x^2 + w_3x^3$

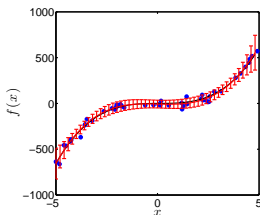
Example

Data sampled from a 3rd order polynomial function:

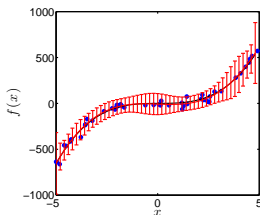
$$t = w_0 + w_1x + w_2x^2 + w_3x^3 + \epsilon$$



Linear



Cubic



6th order

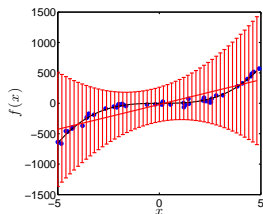
Plots show $t_{\text{new}} \pm \text{var}\{t_{\text{new}}\}$. (Black line is truth).

$$w_0 + \dots + w_6 x^6$$

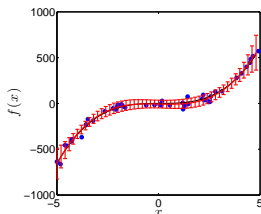
Example

Data sampled from a 3rd order polynomial function:

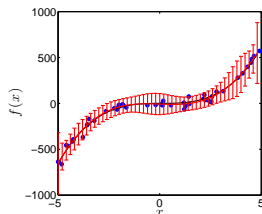
$$t = w_0 + w_1x + w_2x^2 + w_3x^3 + \epsilon$$



Linear



Cubic



6th order

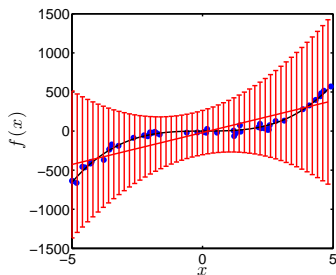
Plots show $t_{\text{new}} \pm \text{var}\{t_{\text{new}}\}$. (Black line is truth).

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

Not complex enough model – more ‘noise’

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

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

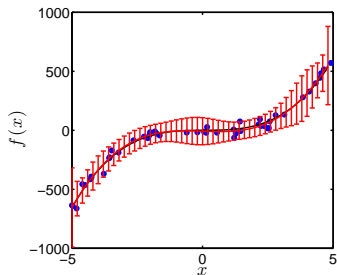


- ▶ The model is too simple.
- ▶ Some true variability can only be modelled noise.
- ▶ $\widehat{\sigma^2}$ is significantly over-estimated.
- ▶ Results in high $\text{var}\{t_{\text{new}}\}$.

Too complex model – parameters not well defined

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

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



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

Olympic prediction

Linear model:

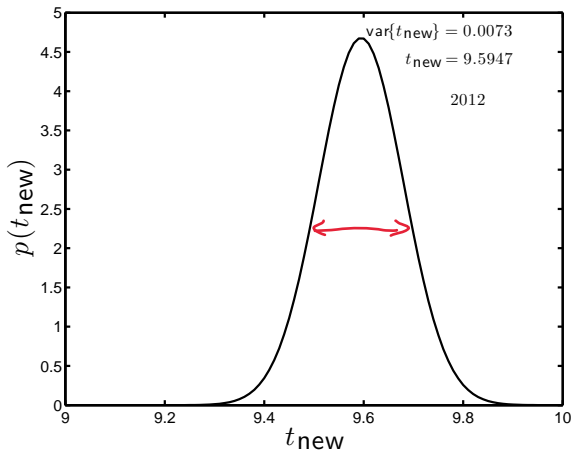
$$t = w_0 + w_1x + \epsilon$$

Olympic prediction

Linear model:

$$\hat{t} = \hat{w}^T x_{new}$$

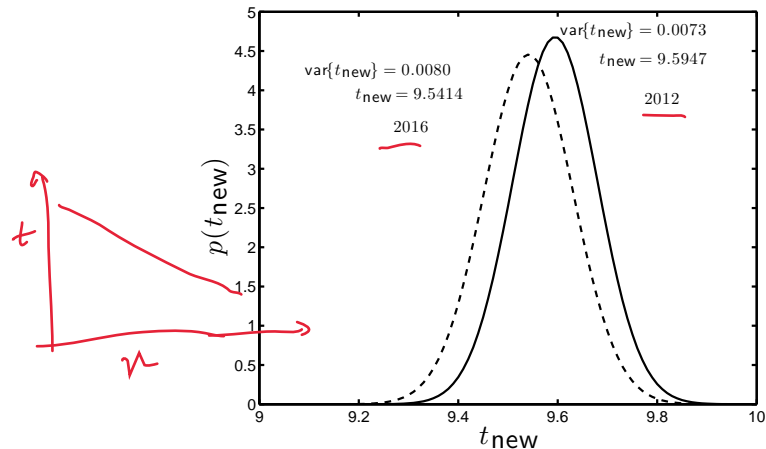
→ $t = w_0 + w_1 x + \epsilon$



Olympic prediction

Linear model:

$$t = w_0 + w_1x + \epsilon$$

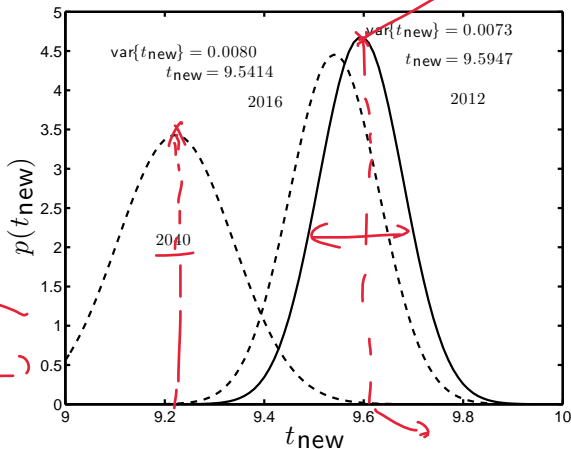


Predictive variance increases as we get further from the training data.

Olympic prediction

Linear model:

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



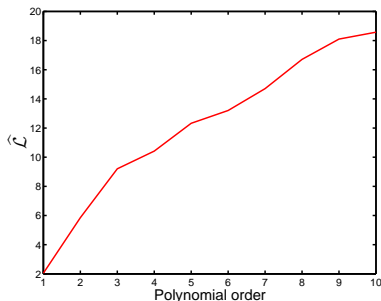
Predictive variance increases as we get further from the training data.

Can we use likelihood to choose models?

- ▶ We've already seen that training loss is no good for model choice.
- ▶ Described cross-validation as an alternative.
- ▶ Can we use the likelihood L or $\log L$?

Can we use likelihood to choose models?

- ▶ We've already seen that training loss is no good for model choice.
- ▶ Described cross-validation as an alternative.
- ▶ Can we use the likelihood L or $\log L$?

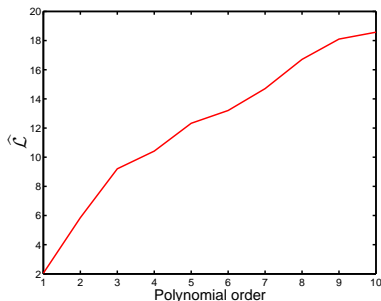


Data from 3rd order polynomial.

- ▶ No.

Can we use likelihood to choose models?

- ▶ We've already seen that training loss is no good for model choice.
- ▶ Described cross-validation as an alternative.
- ▶ Can we use the likelihood L or $\log L$?



Data from 3rd order polynomial.

- ▶ No.
 - ▶ More complex models can always get closer to the 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?

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....

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.

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 random variable (RV):

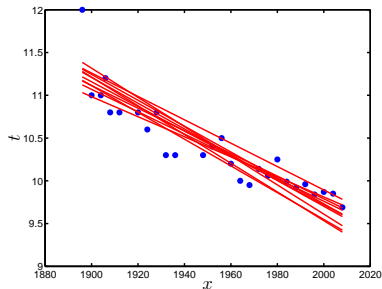
$$p(\mathbf{q}) = \mathcal{N}(\widehat{\mathbf{w}}, \text{cov}\{\widehat{\mathbf{w}}\})$$

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

- ▶ Sample lots of \mathbf{q} from:

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

- ▶ Each corresponds to a model.

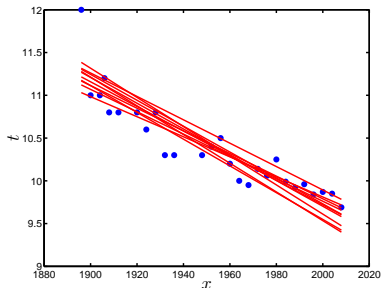


- ▶ Sample lots of \mathbf{q} from:

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

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

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



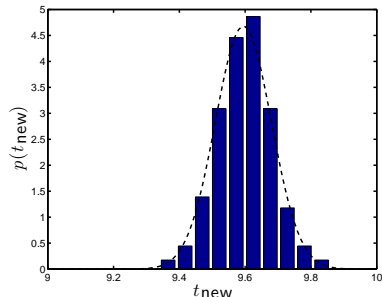
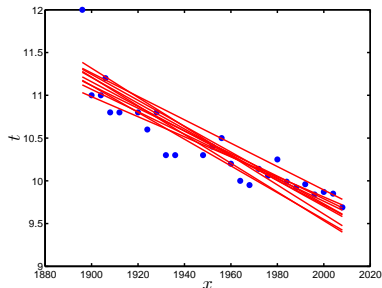
- ▶ Sample lots of \mathbf{q} from:

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

- ▶ 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....