Probabilistic Approach to Linear Regression

Morteza H. Chehreghani morteza.chehreghani@chalmers.se

Department of Computer Science and Engineering Chalmers University of Technology

April 1, 2019

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

Some data and a problem

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



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Recipe for a linear model

More complex model: $t = w_0 + w_1 x + w_2 x^2 + ... + w_D x^D$

$$\mathbf{x}_{n} = \begin{bmatrix} 1 \\ x_{n} \\ x_{n}^{2} \\ \vdots \\ x_{n}^{D} \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{1}^{1} & x_{1}^{2} & \dots & x_{1}^{D} \\ 1 & x_{2}^{1} & x_{2}^{2} & \dots & x_{2}^{D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N}^{1} & x_{N}^{2} & \dots & x_{N}^{D} \end{bmatrix} \mathbf{t} = \begin{bmatrix} t_{1} \\ t_{n} \\ \vdots \\ t_{N} \end{bmatrix},$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Recipe for a linear model

More complex model: $t = w_0 + w_1 x + w_2 x^2 + ... + w_D x^D$

$$\mathbf{x}_{n} = \begin{bmatrix} 1\\ x_{n}\\ x_{n}^{2}\\ \vdots\\ x_{n}^{D} \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 1 & x_{1}^{1} & x_{1}^{2} & \dots & x_{1}^{D}\\ 1 & x_{2}^{1} & x_{2}^{2} & \dots & x_{2}^{D}\\ \vdots & \vdots & \vdots & \ddots & \vdots\\ 1 & x_{N}^{1} & x_{N}^{2} & \dots & x_{N}^{D} \end{bmatrix} \mathbf{t} = \begin{bmatrix} t_{1}\\ t_{n}\\ \vdots\\ t_{N} \end{bmatrix},$$
$$\mathbf{w} = \begin{bmatrix} w_{0}\\ w_{1}\\ \vdots\\ w_{D} \end{bmatrix}, \quad Model : t_{n} = \mathbf{w}^{\mathsf{T}}\mathbf{x}_{n}, \quad or \quad \mathbf{t} = \mathbf{X}\mathbf{w}$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

What about the errors?

$$t = w_0 + w_1 x = \mathbf{w}^{\mathsf{T}} \mathbf{x}$$
$$t = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \ldots + w_D x_D = \sum_{d=0}^{D} w_d x_d = \mathbf{w}^{\mathsf{T}} \mathbf{x}$$

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



▶ ▲母 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ▲ 臣 ● � � �

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:

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

We **should** model the errors

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



They tell us how confident our predictions should be:

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで



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

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

▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで



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

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

<ロト < 回 > < 回 > < 回 > < 回 > < 三 > 三 三



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

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

<ロト < 回 > < 回 > < 回 > < 回 > < 三 > 三 三

What assumptions can we make about ϵ_n ?

It's different for each n.



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

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

ヘロト 人間 とくほとくほとう

э

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



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

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

(日) (四) (日) (日) (日)

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



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

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

- It's different for each n.
- It's positive and negative.
- There doesn't seem to be any relationship between e at different n.
- Looks very hard to model exactly (if it were, it wouldn't be noise!)

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

Our model:

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

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

 $\triangleright \epsilon_n$ is continuous.

• We need to choose $p(\epsilon)$.

Our model:

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

 $\triangleright \epsilon_n$ is continuous.

• We need to choose $p(\epsilon)$.

► Gaussian:



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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Our model:

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

 $\triangleright \epsilon_n$ is continuous.

• We need to choose $p(\epsilon)$.

Gaussian:



▲口▶▲圖▶▲≣▶▲≣▶ = 差 - 釣A@

Gaussian examples



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

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Generating data



900

Evaluate the density:

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

t is a random variable too!

at t = t_n is called for the Likelihood, i.e., the quantity obtained when evaluating the density.

Evaluate the density:

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

t is a random variable too!

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

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

Evaluate the density:

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

t is a random variable too!

- at t = t_n is called for the Likelihood, i.e., the quantity obtained when evaluating the density.
- The higher the value, the more likely t_n is given the model....

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

....the better the model is.



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○○



・ロト・西・・田・・田・・日・



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○○



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

・ロト・西ト・山田・山田・山口・





・ロト・「四ト・「田下・「田下・(日下



シック 単 (中本) (中本) (日)



・ロト・西ト・ヨト・ヨー シック

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

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

- 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!

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

- 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 w and σ² 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:

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

(ロ)、(型)、(E)、(E)、 E) の(()
▶ For each input-response pair, we have a Gaussian likelihood:

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

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

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

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

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

To combine them all, we want the joint likelihood:

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

Assume that the t_n are independent:

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

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

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

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

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

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.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

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

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

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.

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

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Some re-arranging...

$$p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(t_n - \mathbf{w}^\mathsf{T} \mathbf{x}_n)^2\right\}$$
$$\log L = \log \prod_{n=1}^N p(t_n | \mathbf{w}, \mathbf{x}_n, \sigma^2)$$

Some re-arranging...

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

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

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

$$= \sum_{n=1}^N \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \sum_{n=1}^N \frac{1}{2\sigma^2}(t_n - \mathbf{w}^\mathsf{T}\mathbf{x}_n)^2$$

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Looks familiar!

Some re-arranging...

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

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

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

$$= \sum_{n=1}^{N} \log\left(\frac{1}{\sigma\sqrt{2\pi}}\right) - \sum_{n=1}^{N} \frac{1}{2\sigma^2}(t_n - \mathbf{w}^{\mathsf{T}}\mathbf{x}_n)^2$$

$$= -N \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{n=1}^{N} (t_n - \mathbf{w}^{\mathsf{T}}\mathbf{x}_n)^2$$
Looks familiar! To continue (good exercise):

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

The multi-variate Gaussian

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, \ p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{K/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu})\right\}$$

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

The multi-variate Gaussian

$$\begin{split} \mathbf{y} &= \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, \ p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{(2\pi)^{K/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu})\right\} \end{split}$$

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



The multi-variate Gaussian

$$\begin{split} \mathbf{y} &= \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}, \ p(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{(2\pi)^{K/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu})\right\} \end{split}$$

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



The multi-variate Gaussian A special case:

$$\prod_{n=1}^{N} \mathcal{N}(\mu_n, \sigma^2) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_N \end{bmatrix}, \ \boldsymbol{\Sigma} = \begin{bmatrix} \sigma^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}$$

So, in our model:

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Maximising the multi-variate log-likelihood

▶ Partial derivative w.r.t. w, set to zero and solve:

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

Maximising the multi-variate log-likelihood

▶ Partial derivative w.r.t. **w**, set to zero and solve:

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

This is the same expression we've seen before!

Maximising the multi-variate log-likelihood

▶ Partial derivative w.r.t. **w**, set to zero and solve:

$$\log L = \log \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^{2}\mathbf{I})$$

$$\frac{\partial \log L}{\partial \mathbf{w}} = -\frac{1}{2\sigma^{2}}(2\mathbf{X}^{\mathsf{T}}\mathbf{X}\mathbf{w} - 2\mathbf{X}^{\mathsf{T}}\mathbf{t}) = 0$$

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

This is the same expression we've seen before!
 Same for σ²:

$$\frac{\partial \log L}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (\mathbf{t} - \mathbf{X}\mathbf{w})^{\mathsf{T}} (\mathbf{t} - \mathbf{X}\mathbf{w}) = 0$$
$$\sigma^2 = \frac{1}{N} (\mathbf{t} - \mathbf{X}\mathbf{w})^{\mathsf{T}} (\mathbf{t} - \mathbf{X}\mathbf{w})$$

Optimum parameters

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

• Use this to compute optimum $\widehat{\sigma^2}$ from:

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

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

Optimum parameters

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

• Use this to compute optimum $\widehat{\sigma^2}$ from:

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



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?

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

• Imagine there are **true** parameters, **w** and σ^2 .

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

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

- Imagine there are **true** parameters, **w** and σ^2 .
- How good our our estimates $\widehat{\mathbf{w}}$ and $\widehat{\sigma^2}$?
 - 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?

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

- Imagine there are **true** parameters, **w** and σ^2 .
- How good our our estimates $\widehat{\mathbf{w}}$ and $\widehat{\sigma^2}$?
 - 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?

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで



- Imagine there are **true** parameters, **w** and σ^2 .
- How good our our estimates $\widehat{\mathbf{w}}$ and $\widehat{\sigma^2}$?
 - 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?



To progress we need to understand Expectations

- To progress we need to understand Expectations
- lmagine a random variable X with density p(x)



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



- To progress we need to understand Expectations
- lmagine a random variable X with density p(x)
- We want to work out the average value of X, x̃.
- Generate S samples, x_1, \ldots, x_S



- To progress we need to understand Expectations
- lmagine 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, \ldots, x_S
- Average the samples:





Our sample based approximation to x will get better as we take more samples.

Our sample based approximation to x will get better as we take more samples.

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

 We can also (sometimes) compute it exactly using expectations.

Our sample based approximation to x will get better as we take more samples.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

 We can also (sometimes) compute it exactly using expectations.

• Discrete: $\tilde{x} = \mathbf{E}_{p(x)} \{x\} = \sum_{x} xp(x)$

Example:

• X is outcome of rolling die.
$$P(X = x) = 1/6$$

$$\tilde{x} = \sum_{x} x P(X = x) = 3.5$$

Our sample based approximation to x̃ 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.
$$P(X = x) = 1/6$$

•
$$\tilde{x} = \sum_{x} x P(X = x) = 3.5$$

Example:

X is uniform distributed RV between a and b

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

•
$$\tilde{x} = \int_{x=a}^{x=b} x p(x) \, dx = (b-a)/2$$



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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

► In general:

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

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

Some important things:

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

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

► In general:

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

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

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

► In general:

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

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

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @
Expectations – refresher

► In general:

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

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

For vectors of random variables:

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

Expectations - refresher

► In general:

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

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:

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

• Mean: $\mu = \mathsf{E}_{p(\mathsf{x})} \{\mathsf{x}\}$

Covariance:

$$\begin{aligned} \mathsf{cov}\{x\} &= \mathbf{E}_{\rho(\mathbf{x})}\left\{(\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}}\right\} \\ &= \mathbf{E}_{\rho(\mathbf{x})}\left\{\mathbf{x}\mathbf{x}^{\mathsf{T}}\right\} - \mathbf{E}_{\rho(\mathbf{x})}\left\{\mathbf{x}\right\}\mathbf{E}_{\rho(\mathbf{x})}\left\{\mathbf{x}^{\mathsf{T}}\right\} \end{aligned}$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

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|μ, σ²) = N(μ, σ²) Mean: E_{p(x)} {x} = μ Variance: E_{p(x)} {(x − μ)²} = σ² Multi-variate p(x|μ, σ²) = N(μ, Σ) Mean: E_{p(x)} {x} = μ

• Variance:
$$\mathbf{E}_{\rho(\mathbf{x})}\left\{(\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}}\right\} = \boldsymbol{\Sigma}$$

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Parameter estimates:

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

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = のへで

Parameter estimates:

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





Parameter estimates:

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

True values: w, σ²
Our model:

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Parameter estimates:

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

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

Parameter estimates:

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

True values: \mathbf{w}, σ^2

Our model:

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

What's E_{p(t|X,w,σ²)} {ŵ}?
 What do we expect our parameter estimate to be?



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

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

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

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

$$\mathsf{E}_{\rho(\mathsf{t}|\mathsf{X},\mathsf{w},\sigma^2)}\left\{\widehat{\mathsf{w}}
ight\}$$

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

・ロト・(型ト・(型ト・(型ト))

$$\mathsf{E}_{p(\mathsf{t}|\mathsf{X},\mathsf{w},\sigma^2)}\left\{\widehat{\mathsf{w}}
ight\}$$

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

・ロト・(型ト・(型ト・(型ト))

$$\mathsf{E}_{p(\mathsf{t}|\mathsf{X},\mathsf{w},\sigma^2)}\left\{\widehat{\mathsf{w}}
ight\}$$

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

・ロト・(型ト・(型ト・(型ト))

$$\mathsf{E}_{p(\mathsf{t}|\mathsf{X},\mathsf{w},\sigma^2)}\left\{\widehat{\mathsf{w}}\right\}$$

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

$\widehat{\boldsymbol{w}}$ is unbiased

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



• What does $cov{\widehat{\mathbf{w}}}$ tell us?

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

$\operatorname{cov}\{\widehat{\mathbf{w}}\}$

• What does $cov{\widehat{\mathbf{w}}}$ tell us?

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

$\mathrm{cov}\{\widehat{\mathbf{w}}\}$

• What does $cov{\widehat{\mathbf{w}}}$ tell us?

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

$$\operatorname{cov}\{\widehat{\mathbf{w}}\} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ



• What does $cov{\widehat{\mathbf{w}}}$ tell us?

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

$$\operatorname{cov}\{\widehat{\mathbf{w}}\} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

- 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 w₀ and w₁. b how the values should be changed together.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



• What does $cov{\widehat{\mathbf{w}}}$ tell us?

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

$$\operatorname{cov}\{\widehat{\mathbf{w}}\} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

- 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 w₀ and w₁. b how the values should be changed together.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・





▶ What does cov{**ŵ**} tell us?

• Recall the linear model, $\mathbf{w} = \begin{vmatrix} w_0 \\ w_1 \end{vmatrix}$

$$\operatorname{cov}\{\widehat{\mathbf{w}}\} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

- 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 w₀ and w₁. b how the values should be changed together.





$$\begin{array}{lll} \operatorname{cov}\{\widehat{\mathbf{w}}\} & = & \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\mathbf{w}} \widehat{\mathbf{w}}^{\mathsf{T}} \right\} \\ & & - \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\mathbf{w}} \right\} \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\mathbf{w}} \right\}^{\mathsf{T}} \end{array}$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ● ● ●



$$\begin{aligned} \operatorname{cov}\{\widehat{\mathbf{w}}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)}\left\{\widehat{\mathbf{w}}\widehat{\mathbf{w}}^{\mathsf{T}}\right\} \\ &- \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)}\left\{\widehat{\mathbf{w}}\right\} \mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)}\left\{\widehat{\mathbf{w}}\right\}^{\mathsf{T}} \\ &= \mathbf{E}\left\{\widehat{\mathbf{w}}\widehat{\mathbf{w}}^{\mathsf{T}}\right\} - \mathbf{w}\mathbf{w}^{\mathsf{T}} \\ &= \vdots \\ \operatorname{cov}\{\widehat{\mathbf{w}}\} &= \sigma^2(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1} \end{aligned}$$

◆□ ▶ ◆昼 ▶ ◆臣 ▶ ◆臣 ● ● ●

Example



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

 $p(\epsilon_n) = \mathcal{N}(0, \sigma^2)$
 $\sigma^2 = 0.5^2$

Example



▲ロト ▲園 ト ▲ 臣 ト ▲ 臣 ト 一臣 - のへ(で)

Example



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

 $\mathbf{E}_{\rho(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\} - \text{beyond this class} \\ \text{We saw that } \widehat{\mathbf{w}} \text{ was unbiased, what about } \widehat{\sigma^2}?$

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

Useful identity

$$p(\mathbf{t}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$\mathbf{E}_{p(\mathbf{t})} \left\{ \mathbf{t}^{\mathsf{T}} \mathbf{A} \mathbf{t} \right\} = \operatorname{Tr}(\mathbf{A}\boldsymbol{\Sigma}) + \boldsymbol{\mu}^{\mathsf{T}} \mathbf{A} \boldsymbol{\mu}$$
$$\operatorname{Tr}(\mathbf{A}) = \sum_{i} A_{ii}$$

Another useful identity

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

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

Where D is the number of columns in **X** (the number of elements in **w**.

Another useful identity

Tr(AB) = Tr(BA)

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

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

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

► So
$$1 - D/N < 1$$
.

$$\blacktriangleright \text{ So } \widehat{\sigma^2} < \sigma^2$$

• $\widehat{\sigma^2}$ is biased and will generally be too low.

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

 $\blacktriangleright In general D < N.$

$$\blacktriangleright \ \operatorname{So} \ \widehat{\sigma^2} < \sigma^2$$

- $\widehat{\sigma^2}$ 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} .

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

 $\blacktriangleright In general D < N.$

$$\blacktriangleright \ \operatorname{So} \ \widehat{\sigma^2} < \sigma^2$$

- $\hat{\sigma^2}$ 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} .

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ □臣 ○のへ⊙

► As *N* increases,
$$\mathbf{E}_{\rho(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)}\left\{\widehat{\sigma^2}\right\} \rightarrow \sigma^2$$

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

 $\blacktriangleright In general D < N.$

- $\blacktriangleright \ {\rm So} \ \widehat{\sigma^2} < \sigma^2$
- $\hat{\sigma^2}$ is biased and will generally be too low.
- Why?
 - Because it is based on w which will, in general, be closer to the data than w.

- ► As *N* increases, $\mathbf{E}_{\rho(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\} \rightarrow \sigma^2$
- To think about what if D = N or D > N?

Example – beyond this class

Generate 100 datasets from the following model:

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

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

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

Example – beyond this class

Generate 100 datasets from the following model:

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

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



Summary

Computed E_{p(t|X,w,σ²)} {ŵ} = w ŵ is unbiased.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ
- Computed E_{p(t|X,w,σ²)} {ŵ} = w
 ŵ is unbiased.
- Computed $\operatorname{cov}\{\widehat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1}$

Tells us how much slack there is in our parameters.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

- Computed E_{p(t|X,w,σ²)} {ŵ} = w
 ŵ is unbiased.
- Computed $\operatorname{cov}\{\widehat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1}$
 - Tells us how much slack there is in our parameters.

• Computed
$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)} \left\{ \widehat{\sigma^2} \right\} = \sigma^2 (1 - D/N)$$
 [beyond this class!]

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

- $\widehat{\sigma^2}$ is **biased**.
- Gets better and better as we get more data.





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

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで



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



<ロト <回ト < 注ト < 注ト

- 3

Our model is defined as:

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

Given our estimate of the parameters, w and a new input, x_{new}, if we had to predict a single value:

$$t_{\sf new} = \widehat{\mathbf{w}}^{\sf T} \mathbf{x}_{\sf new}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Is this sensible?

Our model is defined as:

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

Given our estimate of the parameters, w and a new input, x_{new}, if we had to predict a single value:

$$t_{new} = \widehat{\mathbf{w}}^{\mathsf{T}} \mathbf{x}_{new}$$

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

$$\mathbf{E}_{\rho(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)}\left\{t_{\mathsf{new}}\right\} = \mathbf{E}_{\rho(\mathbf{t}|\mathbf{X},\mathbf{w},\sigma^2)}\left\{\widehat{\mathbf{w}}^{\mathsf{T}}\mathbf{x}_{\mathsf{new}}\right\} = \mathbf{w}^{\mathsf{T}}\mathbf{x}_{\mathsf{new}}$$

which is a good thing!

・ロト・西ト・山田・山田・山口・

▶ What about var{*t*_{new}}?

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

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = のへで

▶ What about var{*t*_{new}}?

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

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

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

Recall the expression for the covariance of the parameter estimate:

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

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

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

Recall the expression for the covariance of the parameter estimate:

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

Appears in the variance of the prediction:

$$\operatorname{var}\{t_{\operatorname{new}}\} = \mathbf{x}_{\operatorname{new}}^{\mathsf{T}} \operatorname{cov}\{\widehat{\mathbf{w}}\}\mathbf{x}_{\operatorname{new}}\}$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

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

Recall the expression for the covariance of the parameter estimate:

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

Appears in the variance of the prediction:

$$var\{t_{new}\} = \mathbf{x}_{new}^{\mathsf{T}} cov\{\widehat{\mathbf{w}}\}\mathbf{x}_{new}$$

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

Data sampled from a 3rd order polynomial function:

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

Data sampled from a 3rd order polynomial function:

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



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

(日) (四) (日) (日) (日)

Data sampled from a 3rd order polynomial function:

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



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

イロト 不得 トイヨト イヨト

э

Data sampled from a 3rd order polynomial function:

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



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

Data sampled from a 3rd order polynomial function:

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



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

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

・ロト ・ 四ト ・ ヨト ・ ヨト ・ ヨ

Not complex enough model - more 'noise'

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

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



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

Too complex model – parameters not well defined Similarly, we substitute $\widehat{\sigma^2}$ into expression for $\operatorname{cov}\{\widehat{\mathbf{w}}\}$:

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



- 6th order model is too flexible.
- Many sets of parameters lead to a good model.
- Means that cov{ŵ} is high.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Too complex model – parameters not well defined Similarly, we substitute $\widehat{\sigma^2}$ into expression for $\operatorname{cov}\{\widehat{\mathbf{w}}\}$:

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



- 6th order model is too flexible.
- Many sets of parameters lead to a good model.
- Means that cov{ŵ} is high.
 - 'good' 6th order models.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Too complex model – parameters not well defined Similarly, we substitute $\widehat{\sigma^2}$ into expression for $\operatorname{cov}\{\widehat{\mathbf{w}}\}$:

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



- 6th order model is too flexible.
- Many sets of parameters lead to a good model.
- Means that cov{ ŵ } is high.
 - 'good' 6th order models.
 - 'good' 3rd order models.

▶ ▲□ ▶ ▲ 三 ▶ ▲ 三 ▶ ● ○ ○ ○ ○

Linear model:

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

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

Linear model:

 $t = w_0 + w_1 x + \epsilon$



◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

Linear model:

 $t = w_0 + w_1 x + \epsilon$



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

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ○臣 - の々ぐ

Linear model:

 $t = w_0 + w_1 x + \epsilon$



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

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

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?

- 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?





- 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?



No.

More complex models can always get closer to the data.

- 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.

Results in lower $\hat{\sigma}^2$ and higher likelihood.

- Decided to model the noise.
- Recapped random variables.
- Introduced likelihood and maximised it to find $\widehat{\mathbf{w}}$ and $\widehat{\sigma^2}$.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

What did it buy us?

- Decided to model the noise.
- Recapped random variables.
- lntroduced likelihood and maximised it to find $\widehat{\mathbf{w}}$ and σ^2 .

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

- 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....

- Decided to model the noise.
- Recapped random variables.
- lntroduced likelihood and maximised it to find $\widehat{\mathbf{w}}$ and σ^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):

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

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



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

æ

イロト イヨト イヨト

Each corresponds to a model.



Sample lots of **q** from:

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

Each corresponds to a model.

Compute a prediction from each one:

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

ж

イロト イヨト イヨト


Sample lots of **q** from:

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

Each corresponds to a model.

Compute a prediction from each one:

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

Look at the distribution of predictions:



Do we need to take samples at all?

$$\mathbf{E}_{\rho(\mathbf{q})}\left\{t_{\mathsf{new}}\right\} = \int t_{\mathsf{new}} \mathcal{N}(\widehat{\mathbf{w}}, \mathsf{cov}\{\widehat{\mathbf{w}}\}) \ dt_{\mathsf{new}}$$

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ● ●

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