

COM4509/6509

Machine Learning and Adaptive Intelligence

Lecture 4: Linear Regression

Mike Smith* and Matt Ellis

*m.t.smith@sheffield.ac.uk

Detour: Inner and Outer Products

Inner Product

Let \underline{a} and \underline{b} be column vectors, of shape $N \times 1$,

$$\begin{aligned} \mathbf{a} \cdot \mathbf{b} &= \mathbf{a}^T \mathbf{b} = \\ a_1 b_1 + a_2 b_2 + \dots + a_N b_N \end{aligned}$$

Outer Product

Let \underline{a} and \underline{b} be column vectors of shapes $N \times 1$ and $M \times 1$ respectively

$$\mathbf{a} \otimes \mathbf{b} = \mathbf{a} \mathbf{b}^T = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \dots & a_1 b_M \\ a_2 b_1 & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \vdots \\ a_N b_1 & \dots & \dots & a_N b_N \end{bmatrix}$$

Detour: Inner and Outer Products

Inner Product

Let \underline{a} and \underline{b} be column vectors, of shape $N \times 1$

$$\underline{a} \cdot \underline{b} = \underline{a}^T \underline{b} = \\ a_1 b_1 + a_2 b_2 + \dots + a_N b_N$$

Outer Product

Let \underline{a} and \underline{b} be column vectors, of shapes $N \times 1$ and $M \times 1$ respectively

$$\underline{a} \otimes \underline{b} = \underline{a} \underline{b}^T = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \dots & a_1 b_M \\ a_2 b_1 & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ a_N b_1 & \dots & a_N b_M \end{bmatrix}$$

If $\underline{a} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$, $\underline{b} = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}$ what is the inner product of \underline{a} and \underline{b} ?

If $\underline{a} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$, $\underline{b} = \begin{bmatrix} -2 \\ -1 \\ 0 \\ 1 \end{bmatrix}$ what is the outer product of \underline{a} and \underline{b} ?

ACTIVITY:
Compute these!

Detour: Inner and Outer Products

Inner Product

Let \underline{a} and \underline{b} be column vectors, of shape $N \times 1$

$$\underline{a} \cdot \underline{b} = \underline{a}^T \underline{b} = \\ a_1 b_1 + a_2 b_2 + \dots + a_N b_N$$

Outer Product

Let \underline{a} and \underline{b} be column vectors, of shapes $N \times 1$ and $M \times 1$ respectively

$$\underline{a} \otimes \underline{b} = \underline{a} \underline{b}^T = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \dots & a_1 b_M \\ a_2 b_1 & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ a_N b_1 & \dots & a_N b_M \end{bmatrix}$$

If $\underline{a} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$, $\underline{b} = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}$ what is the inner product of \underline{a} and \underline{b} ?

$$\underline{a}^T \underline{b} = [1 \ 2 \ 3] \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} = 1 + 2 - 3 = 0$$

If $\underline{a} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$, $\underline{b} = \begin{bmatrix} -2 \\ -1 \\ 0 \\ 1 \end{bmatrix}$ what is the outer product of \underline{a} and \underline{b} ?

Detour: Inner and Outer Products

Inner Product

Let \underline{a} and \underline{b} be column vectors, of shape $N \times 1$

$$\underline{a} \cdot \underline{b} = \underline{a}^T \underline{b} = \\ a_1 b_1 + a_2 b_2 + \dots + a_N b_N$$

Outer Product

Let \underline{a} and \underline{b} be column vectors, of shapes $N \times 1$ and $M \times 1$ respectively

$$\underline{a} \otimes \underline{b} = \underline{a} \underline{b}^T = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \dots & a_1 b_M \\ a_2 b_1 & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ a_N b_1 & \dots & a_N b_M \end{bmatrix}$$

If $\underline{a} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$, $\underline{b} = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}$ what is the inner product of \underline{a} and \underline{b} ?

$$\underline{a}^T \underline{b} = [1 \ 2 \ 3] \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} = 1 + 2 - 3 = 0$$

If $\underline{a} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$, $\underline{b} = \begin{bmatrix} -2 \\ -1 \\ 0 \\ 1 \end{bmatrix}$ what is the outer product of \underline{a} and \underline{b} ?

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} [-2 \ -1 \ 0 \ 1] = \begin{bmatrix} -2 & -1 & 0 & 1 \\ -4 & -2 & 0 & 2 \end{bmatrix}$$

Linear Regression

We have some training data:

$$x_1 = 3 \quad y_1 = 2$$

$$x_2 = 5 \quad y_2 = 3$$

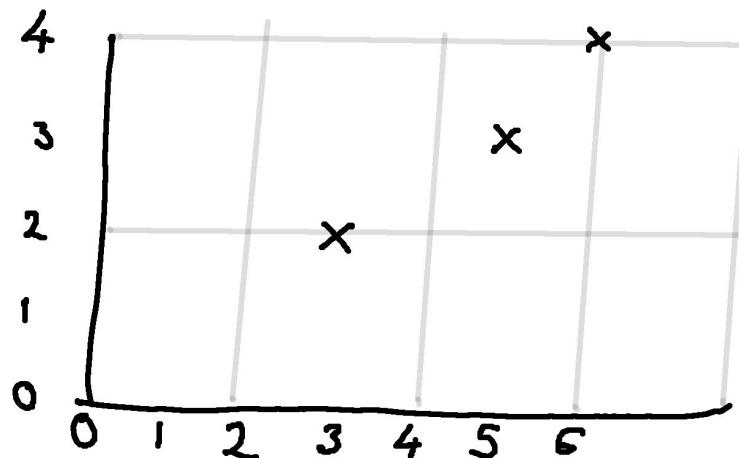
$$x_3 = 6 \quad y_3 = 4$$

We have some training data:

$$x_1 = 3 \quad y_1 = 2$$

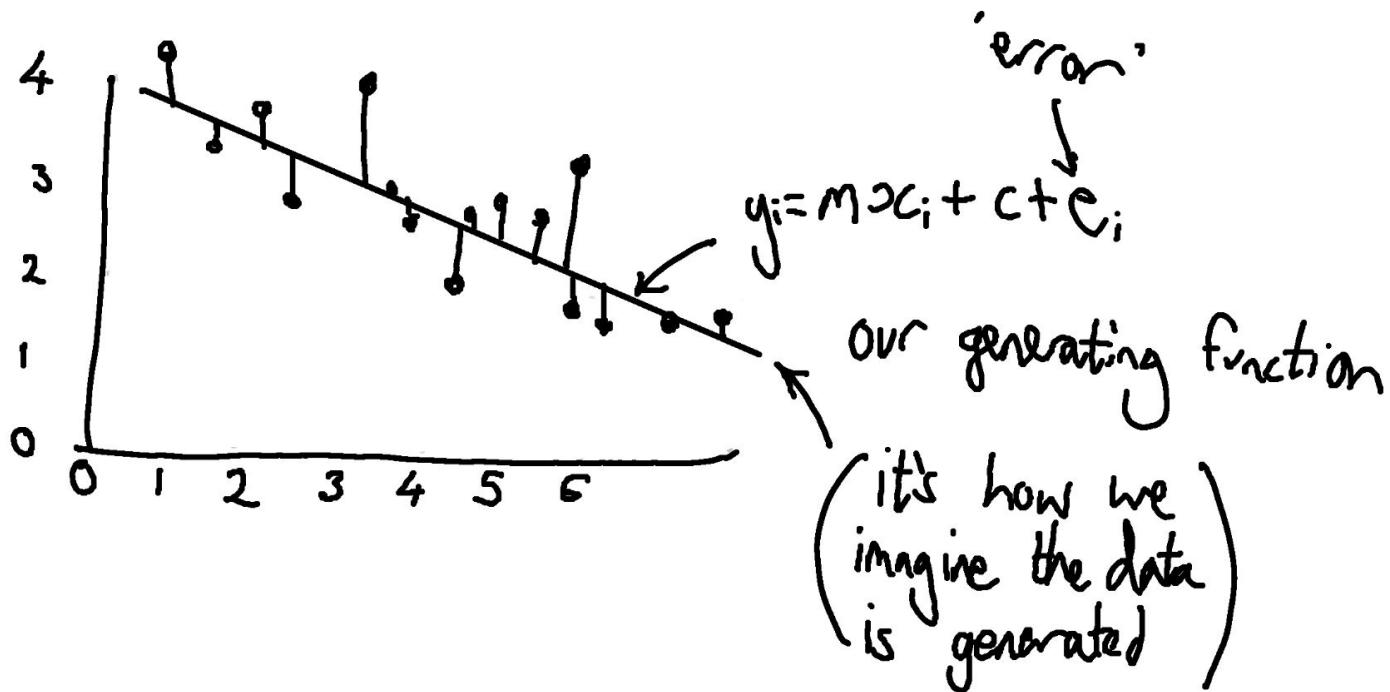
$$x_2 = 5 \quad y_2 = 3$$

$$x_3 = 6 \quad y_3 = 4$$

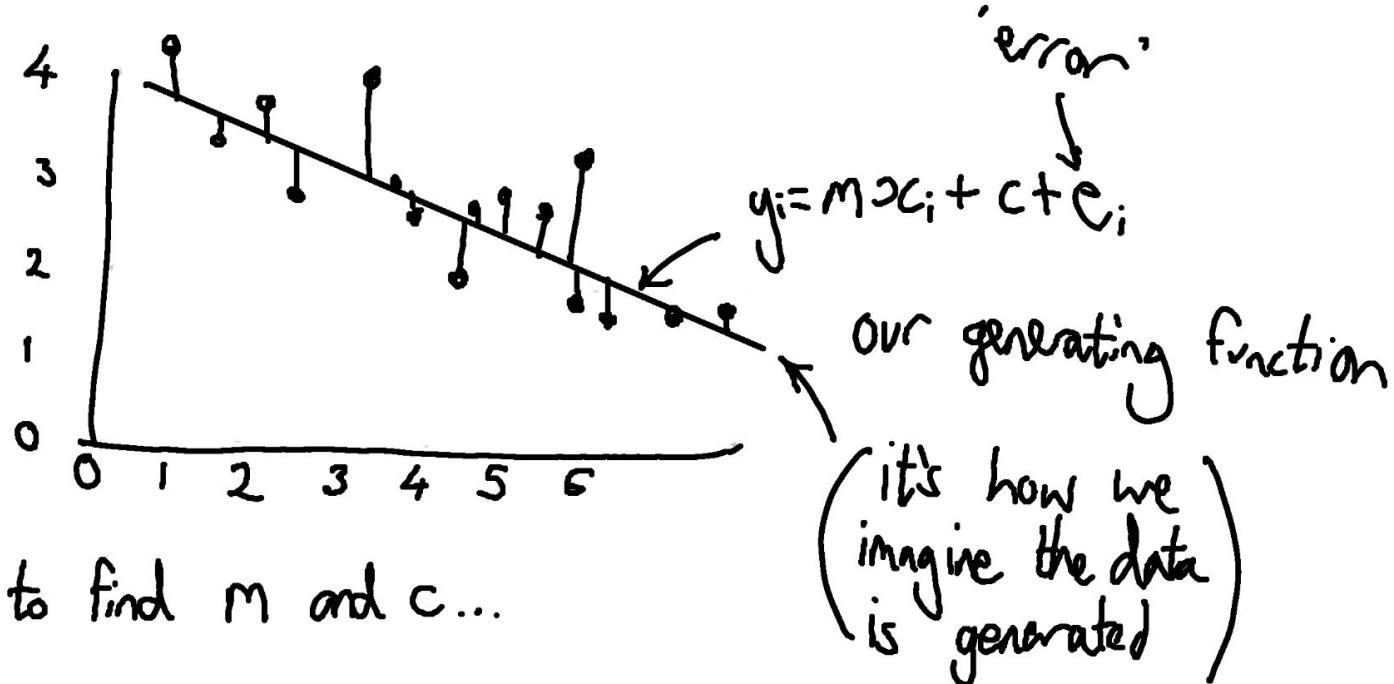


We have a model: Our data is from
a linear, 1st-order polynomial with (Gaussian) noise
added:

We have a model: Our data is from a linear, 1st-order polynomial with (Gaussian) noise added:



We have a model: Our data is from a linear, 1st-order polynomial with (Gaussian) noise added:



We need to find m and c ...

We need a way to say how good (or bad) a choice of m and c is.

We need a way to say how good (or bad) a choice of m and c is.

We will use the sum squared error.

We need a way to say how good (or bad) a choice of m and c is.

We will use the sum squared error. Later we will see why this is a good choice.

We need a way to say how good (or bad) a choice of m and c is.

We will use the sum squared error. Later we will see why this is a good choice.

$$E = \sum_{i=1}^N [(mx_i + c) - y_i]^2$$

We need a way to say how good (or bad) a choice of m and c is.

We will use the sum squared error. Later we will see why this is a good choice.

$$E = \sum_{i=1}^N \left[(mx_i + c) - y_i \right]^2$$

our prediction true label

So this is:

$$E = [(mx_1 + c) - y_1]^2 + [(mx_2 + c) - y_2]^2 + \dots$$

$$E = \sum_{i=1}^N \left[\underbrace{(mx_i + c)}_{\text{our prediction}} - \underbrace{y_i}_{\text{true label}} \right]^2$$

So this is:

$$E = [(mx_1 + c) - y_1]^2 + [(mx_2 + c) - y_2]^2 + \dots$$

This is all quite awkward so it's better
to write with vector notation,

$$E = \sum_{i=1}^N \left[\underbrace{(mx_i + c)}_{\text{our prediction}} - \underbrace{y_i}_{\text{true label}} \right]^2$$

So this is:

$$E = [(mx_1 + c) - y_1]^2 + [(mx_2 + c) - y_2]^2 + \dots$$

This is all quite awkward so it's better to write with vector notation,

First, we look at how to write a vector of:

$$\begin{bmatrix} mx_1 + c \\ mx_2 + c \\ mx_3 + c \\ \vdots \end{bmatrix}$$

$$E = \sum_{i=1}^N \left[\underbrace{(mx_i + c)}_{\text{our prediction}} - \underbrace{y_i}_{\text{true label}} \right]^2$$

If we first make a matrix that looks like this

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \\ \vdots & : \end{bmatrix}$$

If we first make a matrix that looks like this

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \\ \vdots & \vdots \end{bmatrix}$$

then multiply it by a vector $\begin{pmatrix} m \\ c \end{pmatrix}$

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \end{bmatrix} \begin{pmatrix} m \\ c \end{pmatrix} =$$

ACTIVITY: What does this equal?

If we first make a matrix that looks like this

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \\ \vdots & \vdots \end{bmatrix}$$

then multiply it by a vector $\begin{pmatrix} m \\ c \end{pmatrix}$

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \end{bmatrix} \begin{pmatrix} m \\ c \end{pmatrix} = \begin{pmatrix} mx_1 + c \\ mx_2 + c \\ mx_3 + c \end{pmatrix}$$

If we first make a matrix that looks like this

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \\ \vdots & \vdots \end{bmatrix}$$

then multiply it by a vector $\begin{pmatrix} m \\ c \end{pmatrix}$

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \end{bmatrix} \begin{pmatrix} m \\ c \end{pmatrix} = \begin{pmatrix} mx_1 + c \\ mx_2 + c \\ mx_3 + c \end{pmatrix}$$

↑
parameter
vector
w

If we first make a matrix that looks like this

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \\ \vdots & | \end{bmatrix}$$

then multiply it by a vector $\begin{pmatrix} m \\ c \end{pmatrix}$

We call this X ,

The columns are 'regressors' or 'explanatory variables'

often called the design matrix

$$\begin{bmatrix} x_1 & | \\ x_2 & | \\ x_3 & | \end{bmatrix} \begin{pmatrix} m \\ c \end{pmatrix} = \begin{pmatrix} mx_1 + c \\ mx_2 + c \\ mx_3 + c \end{pmatrix}$$

parameter
vector
 $\underline{\omega}$

We write our labels as another vector:

$$\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

We write our labels as another vector:

$$\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

We want to square each row of

$$\underline{y} - \underline{X}\underline{w}$$

We write our labels as another vector:

$$\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

We want to square each row of

$$\underline{y} - \underline{X}\underline{w}$$

We can do this with an inner product:

$$(\underline{y} - \underline{X}\underline{w})^T (\underline{y} - \underline{X}\underline{w})$$

We write our labels as another vector:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

We want to square each row of

$$\mathbf{y} - \mathbf{X}\mathbf{w}$$

We can do this with an inner product:

$$(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

ACTIVITY: Check this is equivalent to the sum of squares error.

$$E = \sum_{i=1}^N \left[\underbrace{(m\mathbf{x}_i + c) - y_i}_{\text{our prediction}} \right]^2 \underbrace{y_i}_{\text{true label}}$$

We write our labels as another vector:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

We want to square each row of

$$\mathbf{y} - \mathbf{X}\mathbf{w}$$

We can do this with an inner product:

$$(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$\begin{bmatrix} y_1 - (m\mathbf{x}_1 + c) & y_2 - m\mathbf{x}_2 + c & \dots \end{bmatrix} \begin{bmatrix} y_1 - (m\mathbf{x}_1 + c) \\ y_2 - (m\mathbf{x}_2 + c) \\ \vdots \end{bmatrix}$$

We write our labels as another vector:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

We want to square each row of

$$\mathbf{y} - \mathbf{X}\mathbf{w}$$

We can do this with an inner product:

$$(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$\begin{bmatrix} y_1 - (m\mathbf{x}_1 + c) & y_2 - m\mathbf{x}_2 + c & \dots \end{bmatrix} \begin{bmatrix} y_1 - (m\mathbf{x}_1 + c) \\ y_2 - (m\mathbf{x}_2 + c) \\ \vdots \end{bmatrix}$$

You can see this leads to the sum of square error we had at the start.

We write our labels as another vector:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

We want to square each row of

$$\mathbf{y} - \mathbf{X}\mathbf{w}$$

We can do this with an inner product:

$$(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\begin{bmatrix} y_1 - (m\mathbf{x}_1 + c) & y_2 - m\mathbf{x}_2 + c & \dots \end{bmatrix} \begin{bmatrix} y_1 - (m\mathbf{x}_1 + c) \\ y_2 - (m\mathbf{x}_2 + c) \\ \vdots \end{bmatrix}$$

You can see this leads to the sum of square error we had at the start.

$$E = \sum_{i=1}^N \left[\underbrace{(m\mathbf{x}_i + c)}_{\text{our prediction}} - \underbrace{y_i}_{\text{true label}} \right]^2$$

We want to minimise our error function,

$$E = (y - X\omega)^T (y - X\omega)$$

We want to minimise our error function,

$$E = (y - X\omega)^T (y - X\omega)$$

with respect to ω (our m & c).

We want to minimise our error function,

$$E = (y - X\omega)^T (y - X\omega)$$

with respect to ω (our m & c).

$$\frac{\partial E}{\partial \omega} = ? = 0$$

We want to minimise our error function,

$$E = (y - X\omega)^T (y - X\omega)$$

with respect to ω (our m & c).

$$\frac{\partial E}{\partial \omega} = ? = 0$$

Check the matrix cookbook...

We want to minimise our error function,

$$E = (y - X\omega)^T (y - X\omega)$$

with respect to ω (our m & c).

$$\frac{\partial E}{\partial \omega} = ? = 0$$

Check the matrix cookbook...

Assume \mathbf{W} is symmetric, then

This identity \longrightarrow
$$\frac{\partial}{\partial s} (\mathbf{x} - \mathbf{As})^T \mathbf{W} (\mathbf{x} - \mathbf{As}) = -2\mathbf{A}^T \mathbf{W} (\mathbf{x} - \mathbf{As}) \quad (84)$$

looks good

We want to minimise our error function,

$$E = (y - X\omega)^T (y - X\omega)$$

with respect to ω (our m & c).

$$\frac{\partial E}{\partial \omega} = ? = 0$$

Check the matrix cookbook...

Assume \mathbf{W} is symmetric, then

This identity \longrightarrow $\frac{\partial}{\partial s} (\mathbf{x} - \mathbf{As})^T \mathbf{W} (\mathbf{x} - \mathbf{As}) = -2\mathbf{A}^T \mathbf{W} (\mathbf{x} - \mathbf{As})$ (84)

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^T \mathbf{x}$	\mathbf{x}
$\mathbf{x}^T \mathbf{w}$	\mathbf{x}
$\mathbf{w}^T \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^T \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also work out from
these rules

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also work out from these rules

(and other rules we have learnt about transposes etc)

$$\frac{\partial}{\partial \underline{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})^\top (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})$$

ACTIVITY: Have a go at finding the value of \mathbf{w} that minimises this error function, E ! We need to differentiate E wrt \mathbf{w} , then set to zero.

Here are some hints:

1. Move the transpose inside the bracket, e.g. by noting that $(\mathbf{X}\mathbf{w})^\top = \mathbf{w}^\top \mathbf{X}^\top$.
2. Multiply out the brackets.
3. Differentiate each term wrt \mathbf{w} .
4. Set equal to zero.
5. Multiply both sides by $(\mathbf{X}\mathbf{X}^\top)^{-1}$ and cancel any constants on both sides etc.

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also
work out from
these rules

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) = (\mathbf{y}^\top - \mathbf{w}^\top \mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\mathbf{w})$$

(moving the transpose inside the bracket).

Move the transpose inside the bracket, e.g. by noting that $(\mathbf{X}\mathbf{w})^\top = \mathbf{w}^\top \mathbf{X}^\top$.

Multiply out the brackets.

Differentiate each term wrt \mathbf{w} .

Set equal to zero.

Multiply both sides by $(\mathbf{X}\mathbf{X}^\top)^{-1}$ and cancel any constants on both sides etc.

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also
work out from
these rules

$$\begin{aligned} \frac{\partial}{\partial \underline{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})^\top (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}}) &= (\mathbf{y}^\top - \underline{\mathbf{w}}^\top \mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\underline{\mathbf{w}}) \\ &= \frac{\partial}{\partial \underline{\mathbf{w}}} (\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\underline{\mathbf{w}} - \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{y} + \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{X}\underline{\mathbf{w}}) \end{aligned}$$

(multiplying
out the
brackets)

Multiply out the brackets.

Differentiate each term wrt \mathbf{w} .

Set equal to zero.

Multiply both sides by $(\mathbf{X}\mathbf{X}^\top)^{-1}$ and
cancel any constants on both
sides etc.

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also work out from these rules

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) = (\mathbf{y}^\top - \mathbf{w}^\top \mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$= \frac{\partial}{\partial \mathbf{w}} \left(\underbrace{\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\mathbf{w}}_0 - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \underbrace{\mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w}}_{2\mathbf{X}^\top \mathbf{X}\mathbf{w}} \right)$$

Differentiate each term wrt \mathbf{w} .

Set equal to zero.

Multiply both sides by $(\mathbf{X}\mathbf{X}^\top)^{-1}$ and
cancel any constants on both
sides etc.

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also work out from these rules

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) = (\mathbf{y}^\top - \mathbf{w}^\top \mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$= \frac{\partial}{\partial \mathbf{w}} \left(\underbrace{\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \mathbf{w}}_0 - \underbrace{\mathbf{w}^\top \mathbf{X}^\top \mathbf{y}}_{-\mathbf{X}^\top \mathbf{y}} + \underbrace{\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}_{2\mathbf{X}^\top \mathbf{X} \mathbf{w}} \right)$$

$$= 2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \mathbf{w} = 0$$

Set equal to zero to find minimum

Set equal to zero.

Multiply both sides by $(\mathbf{X}\mathbf{X}^\top)^{-1}$ and cancel any constants on both sides etc.

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also
work out from
these rules

$$\frac{\partial}{\partial \underline{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})^\top (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}}) = (\mathbf{y}^\top - \underline{\mathbf{w}}^\top \mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})$$

$$= \frac{\partial}{\partial \underline{\mathbf{w}}} \left(\underbrace{\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \underline{\mathbf{w}}}_{\text{o } (\mathbf{X}^\top \mathbf{y})^\top \underline{\mathbf{w}}} - \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{y} + \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}} \right)$$

$$= -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}} = 0$$

$$\Rightarrow 2\mathbf{X}^\top \mathbf{y} = 2\mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}$$

Multiply both sides by $(\mathbf{X}\mathbf{X}^\top)^{-1}$ and
cancel any constants on both
sides etc.

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also
work out from
these rules

$$\frac{\partial}{\partial \underline{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})^\top (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}}) = (\mathbf{y}^\top - \underline{\mathbf{w}}^\top \mathbf{X}^\top) (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})$$

$$= \frac{\partial}{\partial \underline{\mathbf{w}}} \left(\underbrace{\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \underline{\mathbf{w}}}_{\text{o } (\mathbf{X}^\top \mathbf{y})^\top \underline{\mathbf{w}}} - \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{y} + \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}} \right)$$

$$= -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}} = 0$$

$$\Rightarrow 2\mathbf{X}^\top \mathbf{y} = 2\mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}$$

Multiply both sides by $(\mathbf{X}^\top \mathbf{X})^{-1}$ and
cancel any constants on both
sides etc.

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also
work out from
these rules

$$\frac{\partial}{\partial \underline{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})^\top (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}}) = (\mathbf{y}^\top - \underline{\mathbf{w}}^\top \mathbf{X}^\top) (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})$$

$$= \frac{\partial}{\partial \underline{\mathbf{w}}} \left(\underbrace{\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \underline{\mathbf{w}}}_{\textcircled{o}} - \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{y} + \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}} \right)$$

$\mathbf{y}^\top \mathbf{y}$
 $\mathbf{X}^\top \underline{\mathbf{w}}$
 $\mathbf{X}^\top \mathbf{y}$
 $\underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}$

$$= 2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}} = 0$$

$$\Rightarrow 2\mathbf{X}^\top \mathbf{y} = -2\mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}$$

$(\mathbf{X}^\top \mathbf{X})^{-1}$
 $\mathbf{X}^\top \mathbf{y}$
 $\uparrow \text{cancel}$
 $\underline{\mathbf{w}}$

$f(\mathbf{w})$	$\frac{df(\mathbf{w})}{d\mathbf{w}}$
$\mathbf{w}^\top \mathbf{x}$	\mathbf{x}
$\mathbf{x}^\top \mathbf{w}$	\mathbf{x}
$\mathbf{w}^\top \mathbf{w}$	$2\mathbf{w}$
$\mathbf{w}^\top \mathbf{Cw}$	$2\mathbf{Cw}$.

Could also
work out from
these rules

$$\frac{\partial}{\partial \underline{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})^\top (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}}) = (\mathbf{y}^\top - \underline{\mathbf{w}}^\top \mathbf{X}^\top) (\mathbf{y} - \mathbf{X}\underline{\mathbf{w}})$$

$$= \frac{\partial}{\partial \underline{\mathbf{w}}} \left(\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \underline{\mathbf{w}} - \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{y} + \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}} \right)$$

$\cancel{\mathbf{y}^\top \mathbf{y}}$ $\cancel{- \mathbf{y}^\top \mathbf{X} \underline{\mathbf{w}}}$ $\cancel{- \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{y}}$ $\cancel{+ \underline{\mathbf{w}}^\top \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}}$
 $\downarrow \mathbf{y}^\top$ $\downarrow \mathbf{X}^\top \underline{\mathbf{w}}$

$$= 2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}} = 0$$

$$\Rightarrow 2\mathbf{X}^\top \mathbf{y} = -2\mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \underline{\mathbf{w}}$$

$$(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \underline{\mathbf{w}} \quad \begin{matrix} \text{known as} \\ \text{the normal} \\ \text{equation} \end{matrix}$$

Can be solved if

$\mathbf{X}^\top \mathbf{X}$ is invertible.

$$(X^T X)^{-1} X^T y$$

Iterative Optimisation

For the case of linear regression, we've just seen it's got a **closed form solution**. I.e. we can set the gradient to zero, rearrange, and get the answer.

$$(X^T X)^{-1} X^T y$$

Iterative Optimisation

For the case of linear regression, we've just seen it's got a **closed form solution**. I.e. we can set the gradient to zero, rearrange, and get the answer.

For most problems we can't do this. But we might be able to still compute the **gradient** (and use a **gradient descent** algorithm)

$$(X^T X)^{-1} X^T y$$

Iterative Optimisation

For the case of linear regression, we've just seen it's got a **closed form solution**. I.e. we can set the gradient to zero, rearrange, and get the answer.

For most problems we can't do this. But we might be able to still compute the **gradient** (and use a **gradient descent** algorithm)

Side note:

Even though linear regression has a closed form solution, it might be too tricky to compute. ACTIVITY: What size is $(X^T X)$?

$$(X^T X)^{-1} X^T y$$

Iterative Optimisation

For the case of linear regression, we've just seen it's got a **closed form solution**. I.e. we can set the gradient to zero, rearrange, and get the answer.

For most problems we can't do this. But we might be able to still compute the **gradient** (and use a **gradient descent** algorithm)

Side note:

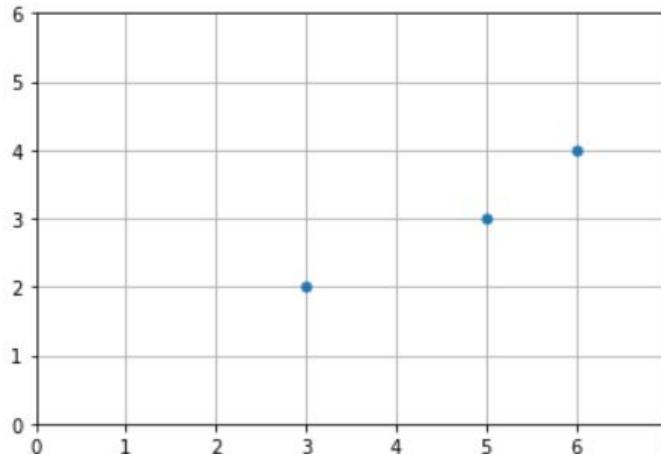
Even though linear regression has a closed form solution, it might be too tricky to compute. **ACTIVITY: What size is $(X^T X)$?**

X is $(N \times D)$, so $X^T X$ is $(D \times D)$.

It takes roughly $O(D^3)$ time to compute the inverse, and uses approximately $O(D^2)$ memory. If our data is high-dimensional (i.e. D is big) then both of these complexities can be a problem.

So we might still need to use gradient descent...

```
1 x = np.array([3,5,6])
2 y = np.array([2,3,4])
3 plt.plot(x,y, '.', markersize=10)
4 plt.xlim([0,7])
5 plt.ylim([0,6])
6 plt.grid()
```



Let's code an example!

```
1 def getprediction(w,x):  
2     X = np.c_[x,np.ones_like(x)]  
3     return X@w
```

This is the method that computes a **prediction** for a given value of **w** and **x**.

```
1 def getprediction(w,x):  
2     X = np.c_[x,np.ones_like(x)]  
3     return X@w
```

Build the design matrix X from our vector x . This has two columns (two regressors). The first is full of x s, the second full of 1s.

This is the method that computes a **prediction** for a given value of w and x .

```
1 def getprediction(w,x):  
2     X = np.c_[x,np.ones_like(x)]  
3     return X@w
```

Our prediction is Xw . In numpy matrix multiplication is with the @ symbol.

Build the design matrix X from our vector x . This has two columns (two regressors). The first is full of x s, the second full of 1s.

This is the method that computes a **prediction** for a given value of w and x .

```
1 def getprediction(w,x):  
2     X = np.c_[x,np.ones_like(x)]  
3     return X@w
```

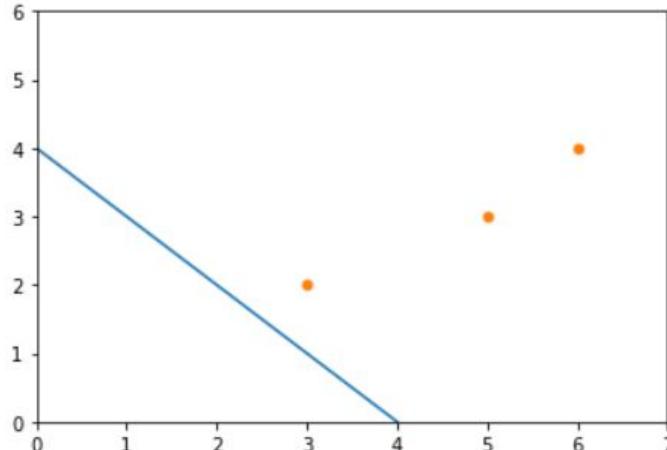
```
1 w = [-1,4]  
2 xtest = np.linspace(0,7,100)  
3 ytest = getprediction(w,xtest)  
4 plt.plot(xtest,ytest, '-')  
5 plt.plot(x,y, '.', markersize=10)  
6 plt.xlim([0,7])  
7 plt.ylim([0,6])
```

Let's plot the predictions
for different values of x ,
with w equal to $[-1,4]^T$.

```
1 def getprediction(w,x):  
2     X = np.c_[x,np.ones_like(x)]  
3     return X@w
```

```
1 w = [-1,4]  
2 xtest = np.linspace(0,7,100)  
3 ytest = getprediction(w,xtest)  
4 plt.plot(xtest,ytest, '-')  
5 plt.plot(x,y, '.', markersize=10)  
6 plt.xlim([0,7])  
7 plt.ylim([0,6])
```

(0.0, 6.0)



Let's plot the predictions
for different values of x ,
with w equal to $[-1,4]^T$.

Not a very
good fit!

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

Let's write the cost function. This returns the sum squared error.

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

```
1 costfn([-1,4],x,y)
```

53

For $w = [-1,4]^T$, $E = 53$.

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

```
1 costfn([-1,4],x,y)
```

53

```
1 costfn([0.5,4],x,y)
```

33.5

For $w = [0.5, 4]^T$, $E = 33.5$.

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

```
1 costfn([-1,4],x,y)
```

53

```
1 costfn([0.5,4],x,y)
```

33.5

For $w = [0.5, 4]^T$, $E = 33.5$.

Ooh, the error is less, it seemed like increasing the first term (the gradient) from -1 to +0.5 helped improved our predictions.

If only there was a way of finding out how much E changes depending on w ...

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

```
1 costfn([-1,4],x,y)
```

53

```
1 costfn([0.5,4],x,y)
```

33.5

```
1 def costfn_grad(w,x,y):  
2     X = np.c_[x,np.ones_like(x)]  
3     return -2*X.T@y + 2*X.T@w
```

This says $-2X^T y + 2X^T Xw$
Which is what we found the
gradient (dE/dw) to be:

$$= -2X^T y + 2X^T Xw$$

Ah! Here's the gradient of
 E wrt w ! We worked this
out in the last section!

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

```
1 costfn([-1,4],x,y)
```

53

```
1 costfn([0.5,4],x,y)
```

33.5

```
1 def costfn_grad(w,x,y):  
2     X = np.c_[x,np.ones_like(x)]  
3     return -2*X.T@y + 2*X.T@w
```

Ah! Here's the gradient of E wrt **w**! We worked this out in the last section!

This says $-2X^T y + 2X^T X w$
Which is what we found the gradient (dE/dw) to be:

$$= -2X^T y + 2X^T X w$$

For complicated functions we use **autodiff** frameworks (TensorFlow, PyTorch, Keras, Theano, JAX, etc). E.g. if the getprediction implemented a deep neural network.

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

```
1 costfn([-1,4],x,y)
```

53

```
1 costfn([0.5,4],x,y)
```

33.5

```
1 def costfn_grad(w,x,y):  
2     X = np.c_[x,np.ones_like(x)]  
3     return -2*X.T@y + 2*X.T@X@w
```

```
1 def costfn_numerical_grad(w,xtest,ytest):  
2     grad = np.array([0,0])  
3     grad[0] = (costfn(w+np.array([0.001,0]),xtest,ytest)-costfn(w-np.array([0.001,0]),xtest,ytest))/0.002  
4     grad[1] = (costfn(w+np.array([0,0.001]),xtest,ytest)-costfn(w-np.array([0,0.001]),xtest,ytest))/0.002  
5     return grad
```

Off topic: To check I've coded my gradient correctly, I'm computing it numerically too...

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

```
1 costfn([-1,4],x,y)
```

53

```
1 costfn([0.5,4],x,y)
```

33.5

```
1 def costfn_grad(w,x,y):  
2     X = np.c_[x,np.ones_like(x)]  
3     return -2*X.T@y + 2*X.T@X@w
```

```
1 def costfn_numerical_grad(w,xtest,ytest):  
2     grad = np.array([0,0])  
3     grad[0] = (costfn(w+np.array([0.001,0]),xtest,ytest)-costfn(w-np.array([0.001,0]),xtest,ytest))/0.002  
4     grad[1] = (costfn(w+np.array([0,0.001]),xtest,ytest)-costfn(w-np.array([0,0.001]),xtest,ytest))/0.002  
5     return grad
```

```
1 costfn_grad([0,4],x,y)
```

array([22, 6])

```
1 def costfn(w,xtest,ytest):  
2     return np.sum((getprediction(w,xtest)-ytest)**2)
```

```
1 costfn([-1,4],x,y)
```

53

```
1 costfn([0.5,4],x,y)
```

33.5

```
1 def costfn_grad(w,x,y):  
2     X = np.c_[x,np.ones_like(x)]  
3     return -2*X.T@y + 2*X.T@X@w
```

```
1 def costfn_numerical_grad(w,xtest,ytest):  
2     grad = np.array([0,0])  
3     grad[0] = (costfn(w+np.array([0.001,0]),xtest,ytest)-costfn(w-np.array([0.001,0]),xtest,ytest))/0.002  
4     grad[1] = (costfn(w+np.array([0,0.001]),xtest,ytest)-costfn(w-np.array([0,0.001]),xtest,ytest))/0.002  
5     return grad
```

```
1 costfn_grad([0,4],x,y)
```

array([22, 6])

```
1 costfn_numerical_grad([0,4],x,y)
```

array([22, 6])

Ok, so we've got our gradient, how do we use it to optimise w?

Gradient Descent

The simplest/most obvious approach is
gradient (steepest) descent.

```
1 w = [-1,4] #initial value
2 for it in range(10000):
3     grad = costfn_grad(w,x,y)
4     w = w - grad*0.01
```

Gradient Descent

The simplest/most obvious approach is **gradient (steepest) descent**.

```
1 w = [-1,4] #initial value
2 for it in range(10000):
3     grad = costfn_grad(w,x,y)
4     w = w - grad*0.01
```

Here we simply subtract (a proportion of) the gradient from the current value of w .

We need to do this in small increments defined by the **learning rate** (aka step size), η [eta]. In this case I've chosen $\eta=0.01$.

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \frac{d}{d\mathbf{w}} E(\mathbf{w}) \Big|_{\mathbf{w}=\mathbf{w}_k}$$

“The new value of w (at iteration $k+1$) equals the old value of w (at iteration k) minus eta (the learning rate) times the gradient of the error E wrt w , at the old value of w .”

Gradient Descent

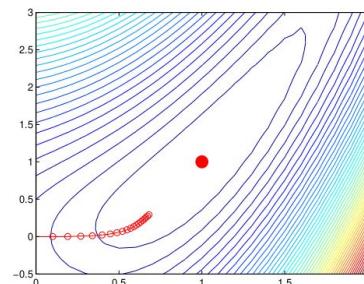
```
1 w = [-1,4] #initial value
2 for it in range(10000):
3     grad = costfn_grad(w,x,y)
4     w = w - grad*0.01
```

Here we simply subtract (a proportion of) the gradient from the current value of w . We need to do this in small increments defined by the **learning rate** (aka step size), η [eta]. In this case I've chosen $\eta=0.01$.

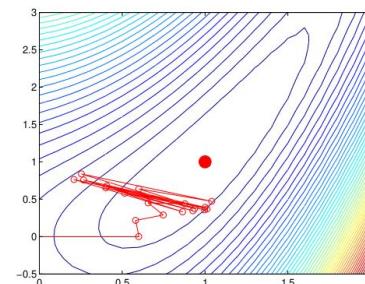
From Mauricio's slide

Step size

- The main issue in gradient descent is how to set the step size.
- If it is too small, convergence will be very slow. If it is too large, the method can fail to converge at all.



(a)



(b)

Figure: The function to optimise is $h(w_1, w_2) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$. The minimum is at (1, 1). In (a) $\eta = 0.1$. In (b) $\eta = 0.6$.

Gradient Descent

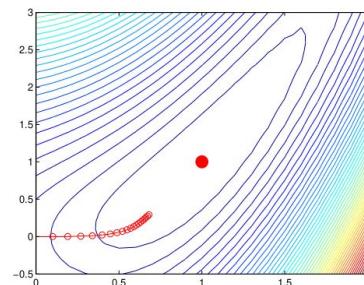
Line Search:

One solution is to pick a direction, and find a point along it where we've reduced the function the most...

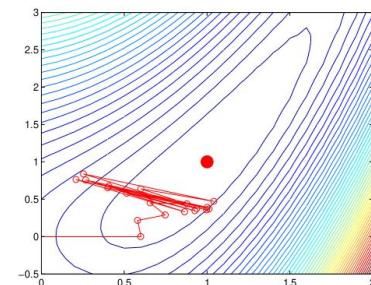
From Mauricio's slide

Step size

- The main issue in gradient descent is how to set the step size.
- If it is too small, convergence will be very slow. If it is too large, the method can fail to converge at all.



(a)



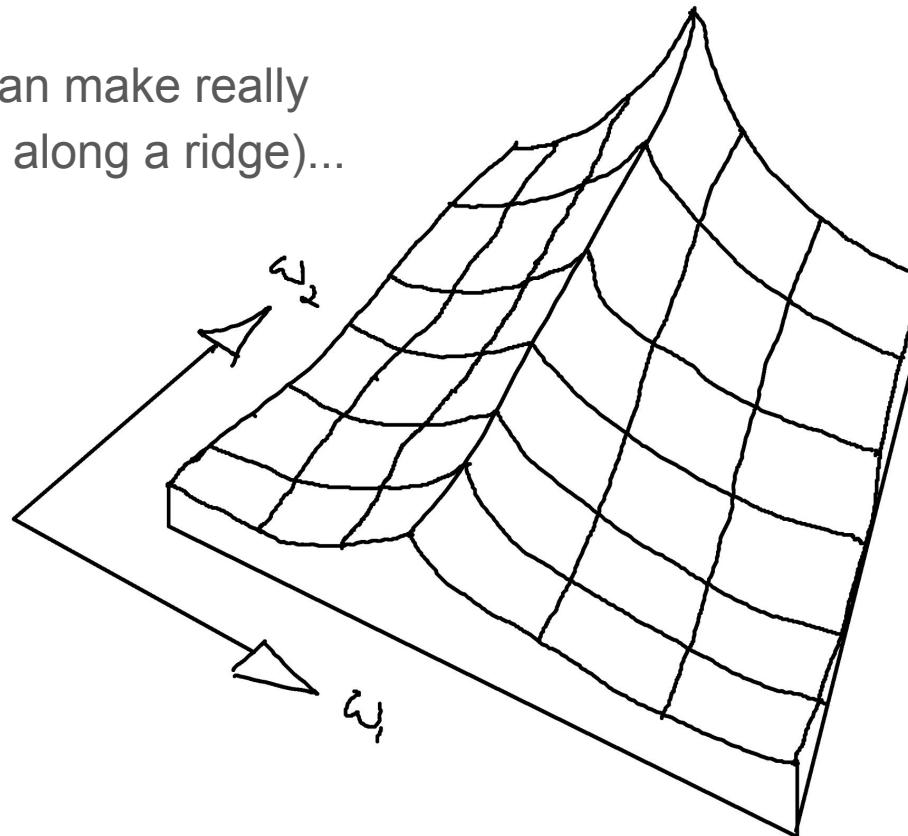
(b)

Figure: The function to optimise is $h(w_1, w_2) = 0.5(w_1^2 - w_2)^2 + 0.5(w_1 - 1)^2$. The minimum is at (1, 1). In (a) $\eta = 0.1$. In (b) $\eta = 0.6$.

Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

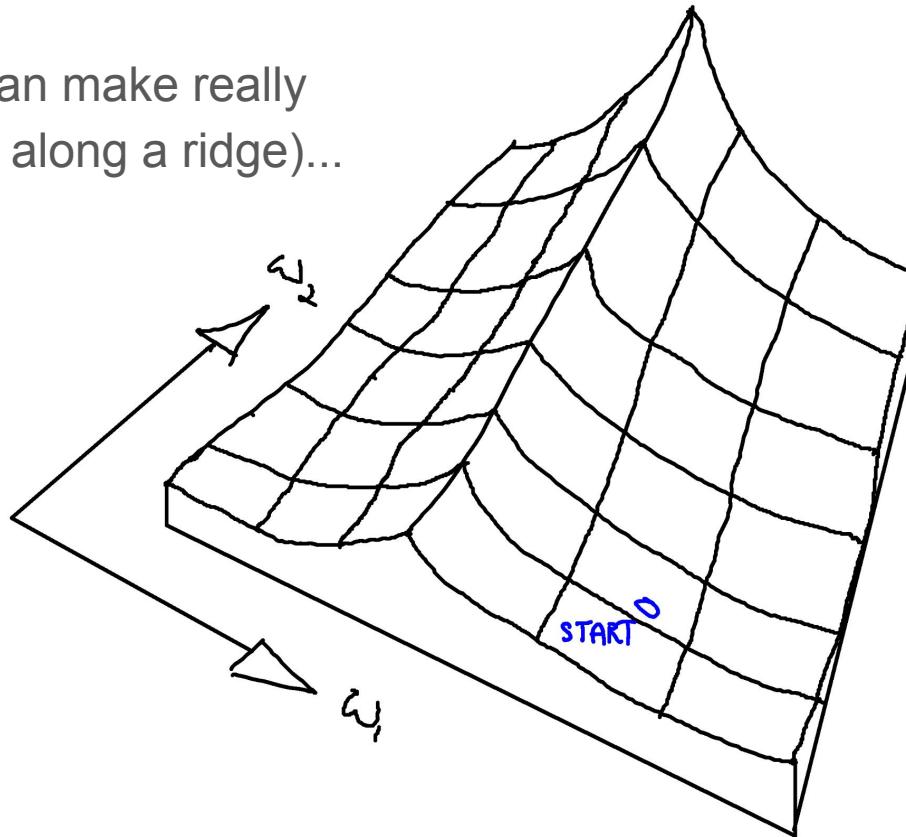
(gradient ascent in this example).



Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

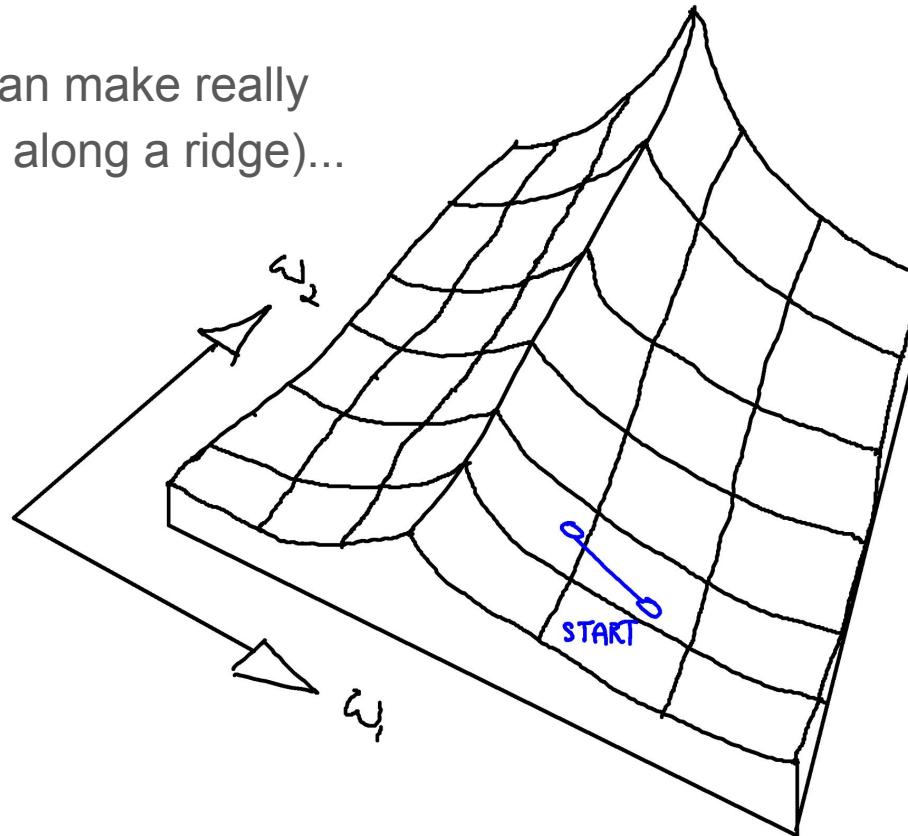
(gradient ascent in this example).



Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

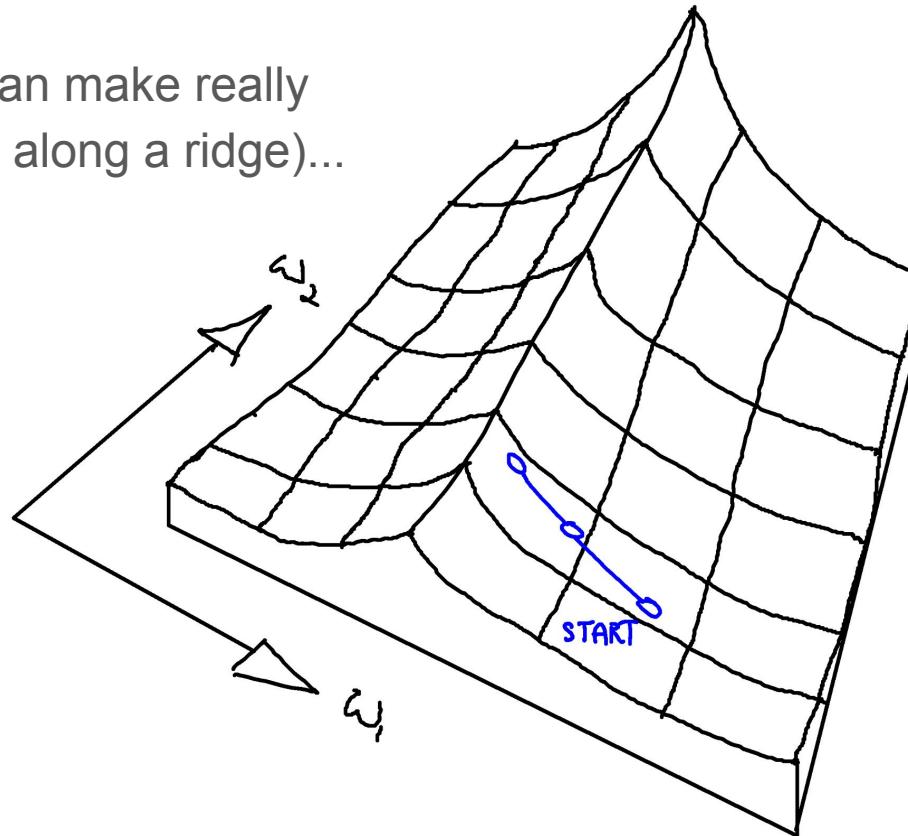
(gradient ascent in this example).



Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

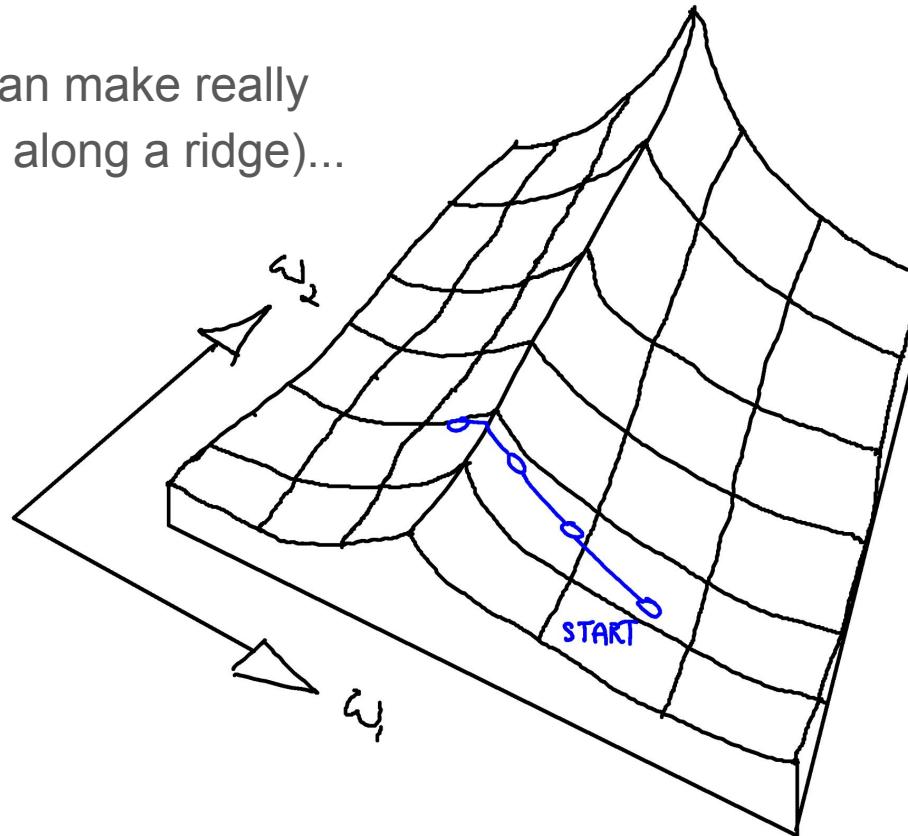
(gradient ascent in this example).



Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

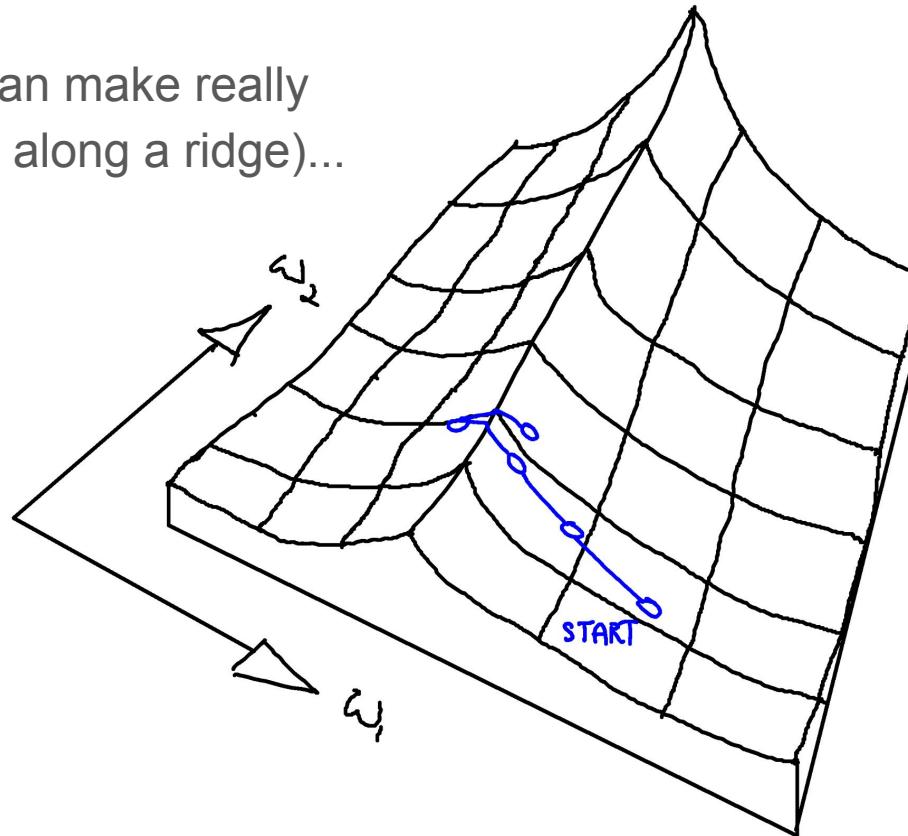
(gradient ascent in this example).



Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

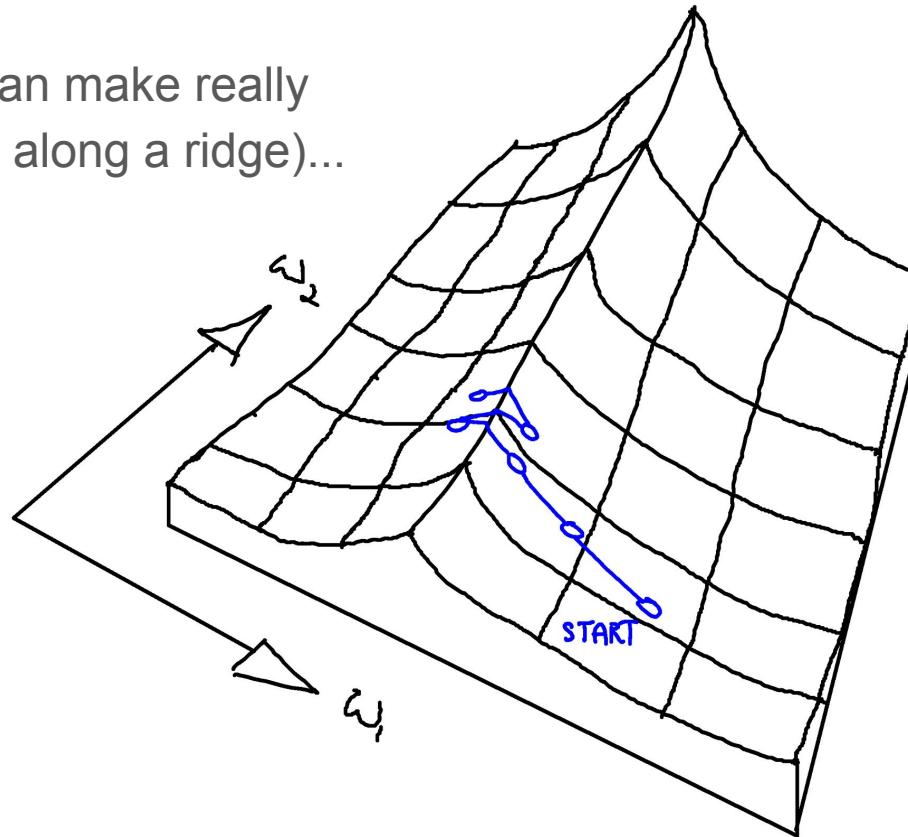
(gradient ascent in this example).



Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

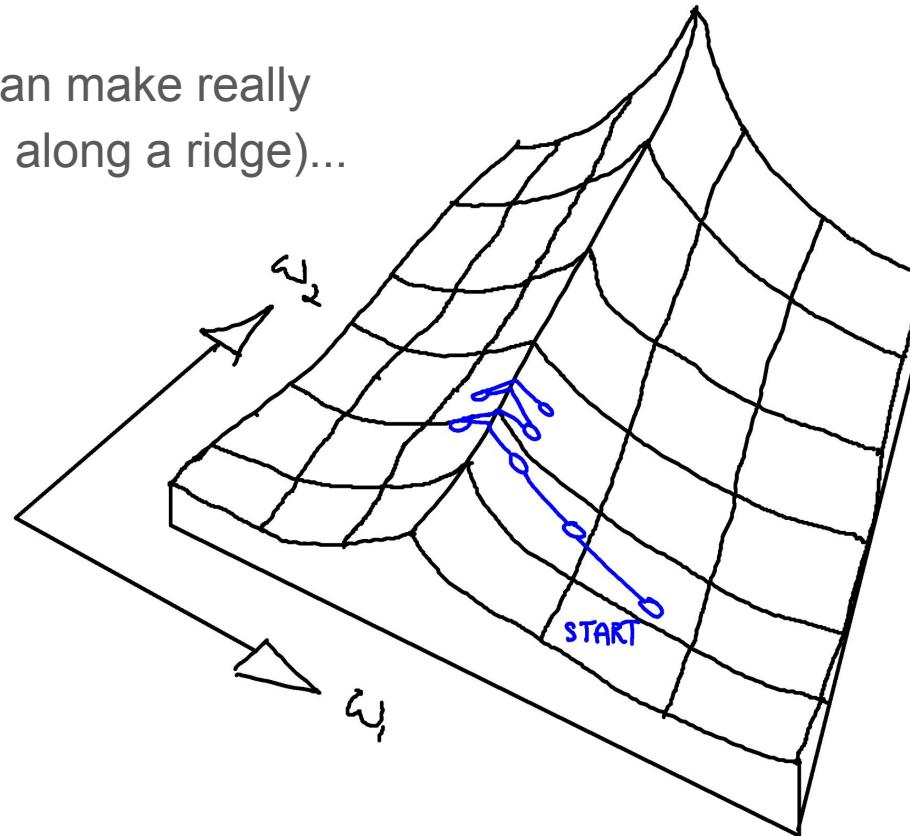
(gradient ascent in this example).



Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

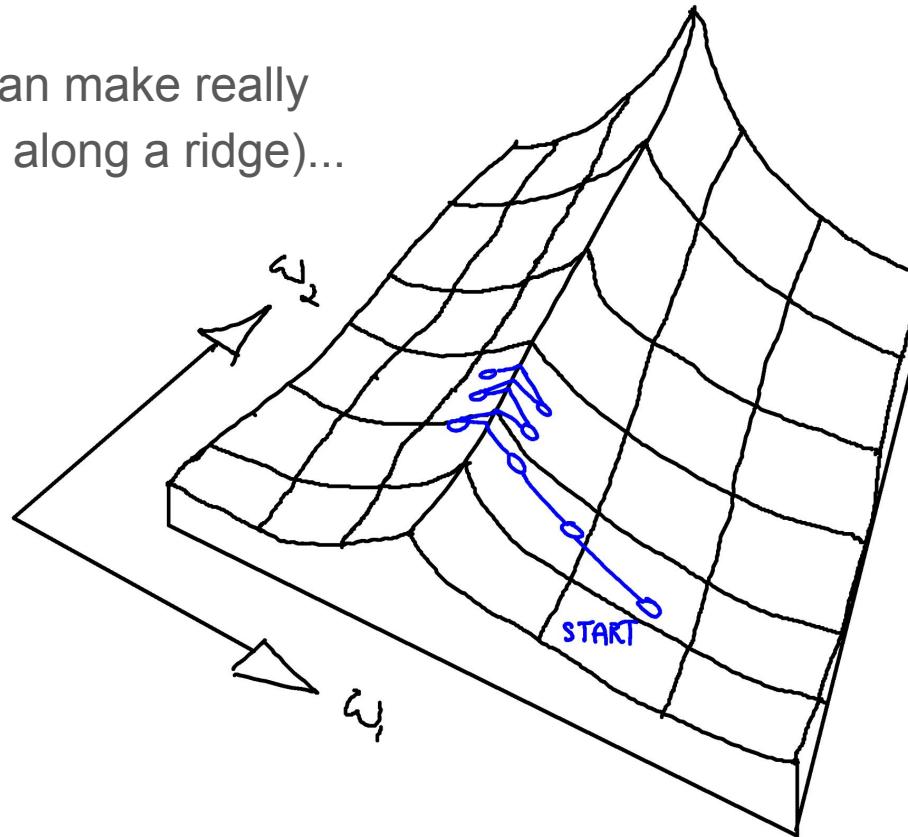
(gradient ascent in this example).



Alternative approaches to Gradient Descent

Gradient descent can make really slow progress (e.g. along a ridge)...

(gradient ascent in this example).



$$x_{k+1} = x_k + t = x_k - \frac{f'(x_k)}{f''(x_k)}$$

Alternative approaches to Gradient Descent

Newton's method [**not examinable!**] will get to a minima quite quickly as it takes into account the 2nd order derivatives:

At each iteration, it amounts to the fitting of a parabola to the graph at the trial value, and then proceeding to the maximum or minimum of that parabola (in higher dimensions, this may also be a saddle point).

– wikipedia

The problem is that this requires finding the Hessian (matrix of 2nd order partial derivatives), which is really expensive.

But approximations to the true Hessian exist. For example the BFGS algorithm.

Alternative approaches to Gradient Descent

- Newton's approach takes too much computation per step
- Steepest descent ends up with many many small steps

An intermediate approach is **conjugate gradient descent**.

Other tricks

- The example gradient descent we computed found the gradient using **all** the data.
 - This is called batch gradient descent.
- If you have a large dataset, or the function you are trying to minimise maybe has a high time/space complexity: consider **mini-batch gradient descent**.
 - This leads to **stochastic gradient descent** (SGD) - as the gradient will depend on the sample chosen at each iteration.

Other tricks

- Choosing the learning rate (in SGD) is now more difficult. Typically some decaying sequence is chosen. If it fulfills the **Robbins-Monro conditions** (and the function we're optimising is differentiable and convex) then, in the limit, it is guaranteed to converge.

$$\eta_k = \frac{1}{k}, \quad \eta_k = \frac{1}{(\tau_0 + k)^\kappa}$$

Example sequences that fulfill the conditions.

$$\sum_{k=1}^{\infty} \eta_k = \infty, \quad \sum_{k=1}^{\infty} \eta_k^2 < \infty.$$

Two key conditions of Robbins-Monro.

Typically though, we use a library optimisation algorithm, such as ADAM, which will adjust the learning rate to help improve convergence.

```

        Ztemp = np.c_[np.tile(Z,[S,1]),np.repeat(np.arange(S),len(Z))]
    if Z.shape[1]==2:
        Ztemp = Z
    self.Z = np.c_[np.tile(Ztemp,[C,1]),np.repeat(np.arange(self.C),len(Ztemp))]

    if likemode=='distribution' or likemode=='process':
        assert gpflowkernellike is not None, "You need to specify the kernel to use a distribution or process"
        self.klike = Kernel(gpflowkernellike)

        self.Zlike = np.c_[Ztemp,np.repeat(0,len(Ztemp))]
        if likelr is None: likelr = lr * 4 #probably can optimise this a little quicker?
        self.likeoptimizer = tf.keras.optimizers.Adam(learning_rate=likelr,amsgrad=False)

    self.fullX = X

```

```

        if (self.mulike is None) or (it%50<25): #optimise latent fns
            gradients = tape.gradient(elbo_loss, [self.mu,self.scale])
            self.optimizer.apply_gradients(zip(gradients, [self.mu, self.scale]))
        else: #this optimises the likelihood...
            if self.likemode=='distribution':
                gradients = tape.gradient(elbo_loss, [self.mulike])
                self.likeoptimizer.apply_gradients(zip(gradients, [self.mulike]))
            if self.likemode=='process':
                gradients = tape.gradient(elbo_loss, [self.mulike,self.scalelike])
                self.likeoptimizer.apply_gradients(zip(gradients, [self.mulike,self.scalelike]))

```

Snippet of a module I wrote for calibrating air pollution sensors.

Typically though, we use a library optimisation algorithm, such as ADAM, which will adjust the learning rate to help improve convergence.

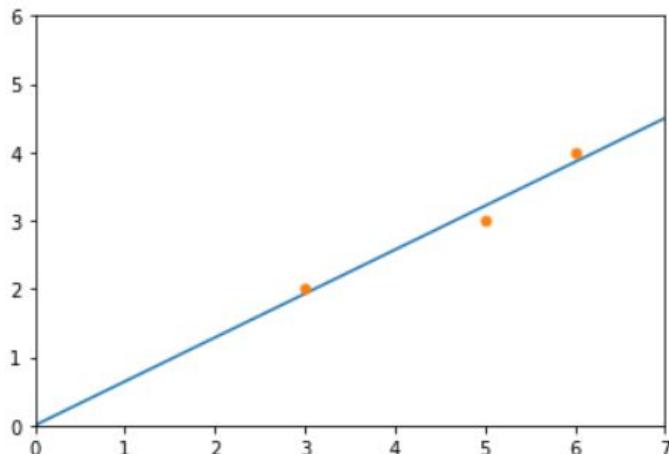
Other tricks

Finally note: Most optimisers will work better if the inputs (features) are normalised.

```
1 w = [-1,4] #initial value
2 for it in range(10000):
3     grad = costfn_grad(w,x,y)
4     w = w - grad*0.01
```

```
1 xtest = np.linspace(0,7,10)
2 ytest = getprediction(w,xtest)
3 plt.plot(xtest,ytest,'-')
4 plt.plot(x,y,'.',markersize=10)
5 plt.xlim([0,7])
6 plt.ylim([0,6])
```

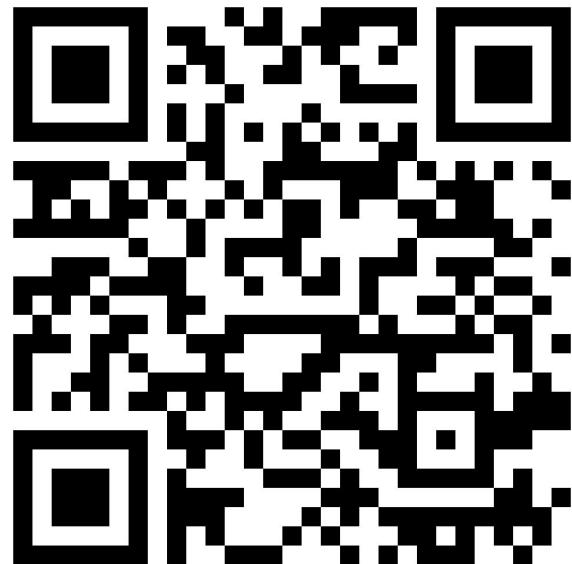
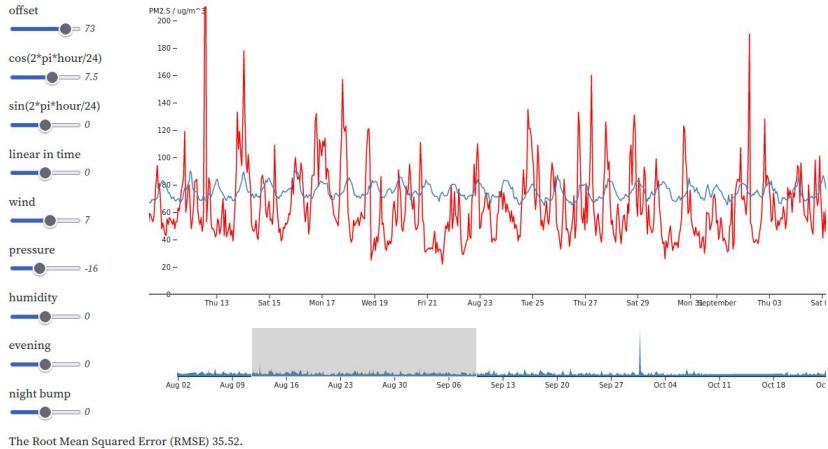
(0.0, 6.0)



break

Linear Regression

The simplest way to make a prediction is with linear regression. We have a list of covariates or inputs that we believe might help predict the pollution (time and date, humidity, etc). We simply take each of these and scale it then add them all together. If we scale them correctly we hopefully will be able to get close to the true pollution.



Linear regression: try predicting the air pollution in [this example](#) by manually adjusting the parameters (weights) that say how much to scale each feature.

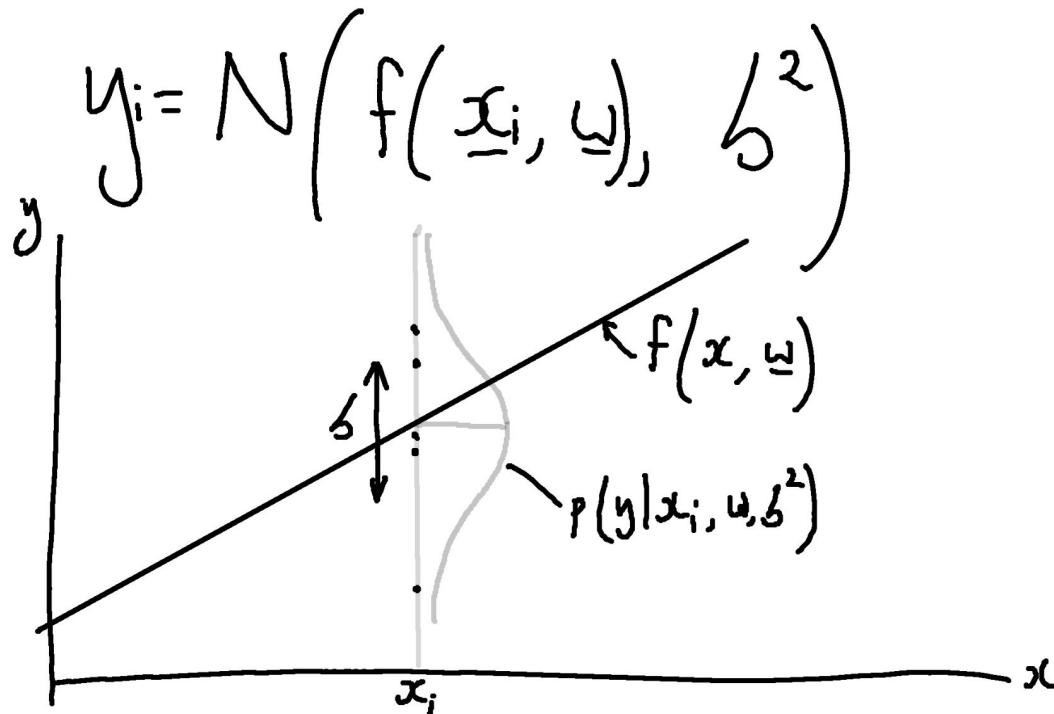
Maximum Likelihood Estimation

We said that our noise (error) term was assumed to be Gaussian distributed.

If we assume this, then we can come at the problem from a different angle: Trying to maximise the probability that the model generated the data, by adjusting model parameters.

We said our noise (error) term
was assumed to be Gaussian distributed,

We said our noise (error) term
was assumed to be Gaussian distributed,



So, for a given x_i and w , we can say
how probable a given value of y_i is:

So, for a given x_i and w , we can say
how probable a given value of y_i is:

$$p(y_i | x_i, w, \sigma^2) = N(y | f(x_i, w), \sigma^2)$$

So, for a given x_i and w , we can say
how probable a given value of y_i is:

$$p(y_i | x_i, w, \sigma^2) = N(y | f(x_i, w), \sigma^2)$$

If we assume our data points are independent, then

$$p(y_1 \dots y_n | x_1 \dots x_n, w, b^2) = \prod_{i=1}^n p(y_i | x_i, w, b^2)$$

So, for a given x_i and w , we can say
how probable a given value of y_i is:

$$p(y_i | x_i, w, \sigma^2) = N(y | f(x_i, w), \sigma^2)$$

If we assume our data points are independent, then

$$p(y_1 \dots y_n | x_1 \dots x_n, w, b^2) = \prod_{i=1}^n p(y_i | x_i, w, b^2)$$

We assume our errors follow the same distribution:

i.i.d. \rightarrow independent
&
identically
distributed

So, for a given x_i and w , we can say
how probable a given value of y_i is:

$$p(y_i | x_i, w, \sigma^2) = N(y | f(x_i, w), \sigma^2)$$

If we assume our data points are independent, then

$$p(y_1 \dots y_n | x_1 \dots x_n, w, b^2) = \prod_{i=1}^n p(y_i | x_i, w, b^2)$$

We assume our errors follow the same distribution:

i.i.d. \rightarrow independent
&
identically
distributed

|
In this case
a Normal distribution
with variance σ^2

So, for a given x_i and w , we can say how probable a given value of y_i is:

$$p(y_i | x_i, w, \sigma^2) = N(y_i | f(x_i, w), \sigma^2)$$

If we assume our data points are independent, then

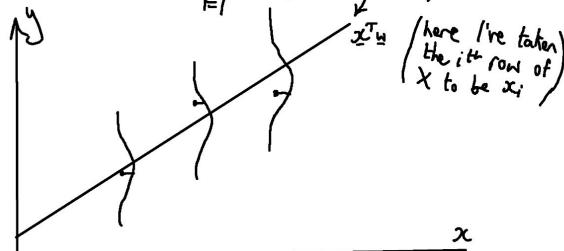
$$p(y_1 \dots y_n | x_1 \dots x_n, w, b^2) = \prod_{i=1}^n p(y_i | x_i, w, b^2)$$

We assume our errors follow the same distribution:

i.i.d. \rightarrow independent & identically distributed

In this case
a Normal distribution
with variance σ^2

$$p(y | X, w, b^2) = \prod_{i=1}^N N(y_i | x_i^T w, b^2)$$



Likelihood

Before, we wanted to minimise a cost function. Here we want to maximise the probability that the data came from our model (by adjusting our model). Because we are interested in this probability as a function of w and σ^2 , it won't integrate to one over these parameters. So it's not a probability distribution. We instead call it the **likelihood function**,

$$L(\underline{w}, b^2) = \prod_{i=1}^N N(y | \underline{x}_i^T \underline{w}, b^2)$$

Maximum likelihood estimation selects parameters (in this case, w and σ^2) to maximise the likelihood function.

PDF for a Normal

$$\frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

The Normal Distribution

Looks tricky, but
let's step
through how
you might work
this out from
scratch...

PDF for a Normal

$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

Looks tricky, but
let's step
through how
you might work
this out from
scratch...

PDF for a Normal

$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Gaussian
shape \rightarrow e^{-x^2}

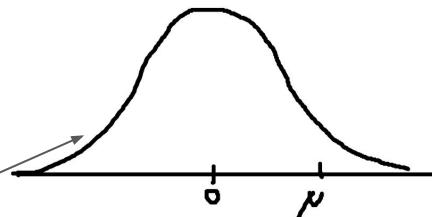
Looks tricky, but
let's step
through how
you might work
this out from
scratch...

PDF for a Normal

$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Gaussian
shape

$$e^{-x^2}$$



Looks tricky, but
let's step
through how
you might work
this out from
scratch...

PDF for a Normal

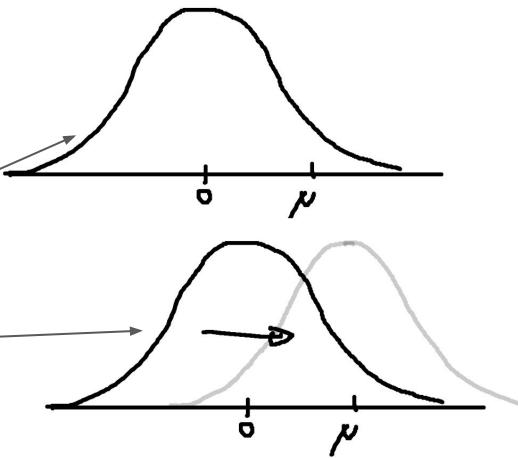
$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Gaussian
shape \rightarrow

shift it
by mean \rightarrow

$$e^{-x^2}$$

$$e^{-(x-\mu)^2}$$



PDF for a Normal

$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Looks tricky, but let's step through how you might work this out from scratch...

Gaussian shape

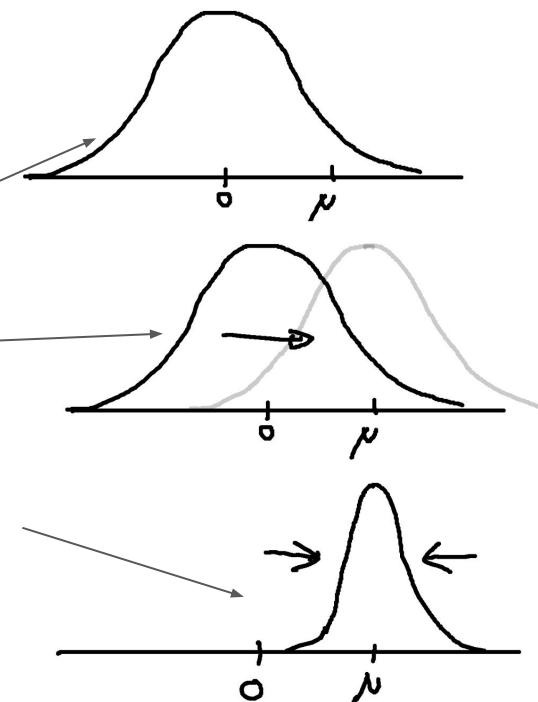
shift it by mean

Make it narrow or wide

$$e^{-x^2}$$

$$e^{-(x-\bar{x})^2}$$

$$e^{-\frac{1}{2}\left(\frac{x-\bar{x}}{\sigma}\right)^2}$$



Looks tricky, but let's step through how you might work this out from scratch...

PDF for a Normal

$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

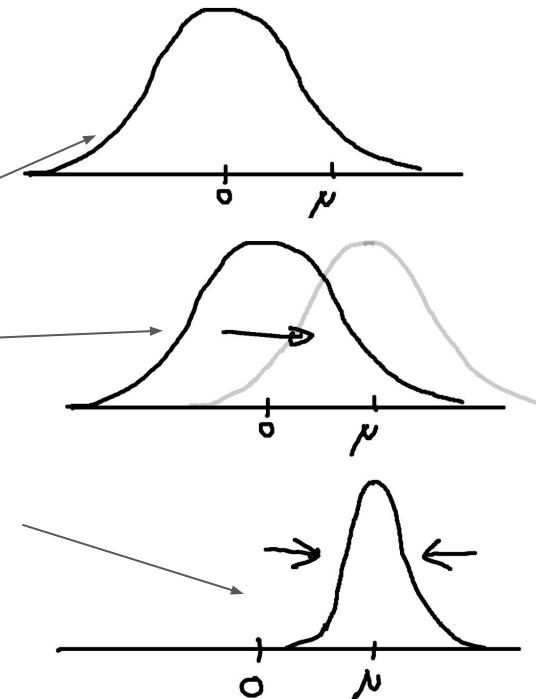
Gaussian shape

$$e^{-x^2}$$

shift it by mean

$$e^{-(x-\mu)^2}$$

$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$



Need a normalising term:
E.g. If σ is big, then our Gaussian is wide, so we need to make it less tall.

Let's go back now to our expression for the likelihood

$$L(\omega, b^2) = \prod_{i=1}^N N(y | x_i^\top \omega, b^2)$$

Remember we had assumed iid for our noise, so we can say that the joint probability of all the observations is the product of the individual probabilities.

Let's go back now to our expression for the likelihood

$$\mathcal{L}(\underline{\omega}, \sigma^2) = \prod_{i=1}^N N(y | \underline{x}_i^\top \underline{\omega}, \sigma^2)$$

Maximum Likelihood Estimation selects
 $\underline{\omega}$ & σ^2 to maximise \mathcal{L} :

Let's go back now to our expression for the likelihood

$$L(\underline{\omega}, \sigma^2) = \prod_{i=1}^N N(y | \underline{x}_i^\top \underline{\omega}, \sigma^2)$$

Maximum Likelihood Estimation selects
 $\underline{\omega}$ & σ^2 to maximise L :

Let's substitute in the expression for
the Normal distribution...

$$L(\underline{\omega}, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^\top \underline{\omega}}{\sigma} \right)^2}$$

Log Likelihood

We want to differentiate the likelihood with respect to w and σ^2 , but in its current form that is quite tricky.

$$L(w, b) = \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - x_i^T w}{\sigma} \right)^2}$$

Log Likelihood

We want to differentiate the likelihood with respect to w and σ^2 , but in its current form that is quite tricky.

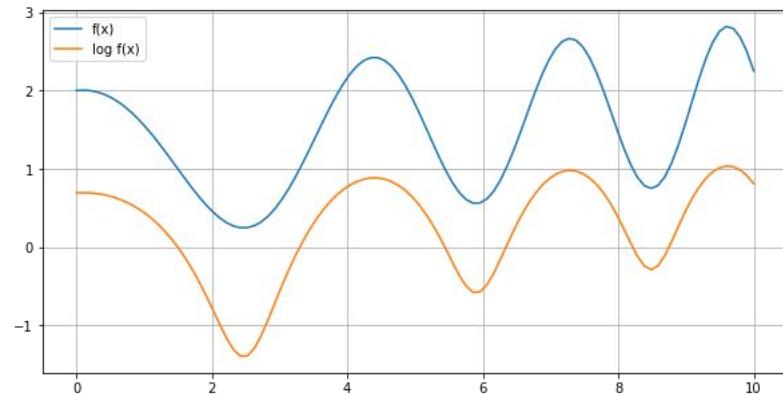
The location of the maximum of L is in the same place as the maximum of $\log L$ (as log is monotonic).

So we first, find the log of the likelihood, then differentiate *that*...

ACTIVITY: Compute the log likelihood.

Hint: Start by thinking about how the log product of things is the sum of the log of things.

$$L(w, b^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{(y_i - x_i^T w)^2}{b^2} \right)}$$



An example function, and its log, to demonstrate that the maxima and minima are in the same places!

The log likelihood (LL)

$$\log L(\underline{\omega}, b^2) = \sum_{i=1}^N \log \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{\omega}^T \underline{x}_i}{b} \right)^2}$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}$$

differentiate...

$$y^T y - y^T X_w - w^T X^T y + w^T X^T X_w$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}$$

differentiate...

$$y^T y - y^T X_w - w^T X^T y + w^T X^T X_w$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}$$

differentiate...

$$y^T y - y^T X_w - w^T X^T y + w^T X^T X_w$$

$$\frac{\partial}{\partial w} = X^T y - X^T y + 2X^T X_w$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}$$

differentiate...

$$\begin{aligned} & y^T y - y^T X_w - w^T X^T y + w^T X^T X_w \\ & \frac{\partial}{\partial w} = X^T y - X^T y + 2 X^T X_w \\ & = -2 X^T (y - X_w) \end{aligned}$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}$$

differentiate...

$$y^T y - y^T X_w - w^T X^T y + w^T X^T X_w$$

$$\frac{\partial}{\partial w} = X^T y - X^T y + 2X^T X_w$$

$$= 2X^T (y - X_w)$$

$$X^T y = X^T X_w$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}$$

differentiate...

$$y^T y - y^T X_w - w^T X^T y + w^T X^T X_w$$

$$\frac{\partial}{\partial w} = X^T y - X^T y + 2X^T X_w$$

$$= 2X^T (y - X_w)$$

$$X^T y = X^T X_w$$

$$(X^T X)^{-1} X^T y = w$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}_{y^T y - y^T X_w - w^T X^T y + w^T X^T X_w}$$

Same as before!

$$\begin{aligned} \frac{\partial}{\partial w} &= X^T y - X^T y + 2X^T X_w \\ &= -2X^T (y - X_w) \\ X^T y &= X^T X_w \end{aligned}$$

$$(X^T X)^{-1} X^T y = w$$

The log likelihood (LL)

$$\log L(\underline{w}, b^2) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= \sum_{i=1}^N -\log \sqrt{2\pi} + \log e^{-\frac{1}{2} \left(\frac{y_i - \underline{x}_i^T \underline{w}}{b} \right)^2}$$

$$= -N \log \sqrt{2\pi} - \frac{1}{2b^2} \sum (y_i - \underline{x}_i^T \underline{w})^2$$

Let's write using vector notation...

$$LL = -N \log \sqrt{2\pi} - \frac{1}{2b^2} \underbrace{(y - X_w)^T (y - X_w)}_{y^T y - y^T X_w - w^T X^T y + w^T X^T X_w}$$

Same as before!

$$\frac{\partial}{\partial w} = X^T y - X^T y + 2X^T X_w$$

$$= -2X^T (y - X_w)$$

$$X^T y = X^T X_w$$

$$(X^T X)^{-1} X^T y = w$$

We can differentiate and solve for σ^2 in a similar way.

Summary: Maximum Likelihood Estimation and Ordinary Least Squares

If we have Gaussian-distributed i.i.d. noise, minimising the sum squared error (i.e. “ordinary least squares”) is the maximum likelihood estimator (for linear regression).

Basis Functions

We've already used a basis, the 1st-order polynomial had two bases: x^1 and x^0 .

If we first make a matrix that looks like this

$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ \vdots & \vdots \end{bmatrix}$$

then multiply it by a vector $\begin{pmatrix} m \\ c \end{pmatrix}$

$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \end{bmatrix} \begin{pmatrix} m \\ c \end{pmatrix} =$$

ACTIVITY: What does this equal?

From slide 16...

Basis Functions

We've already used a basis, the 1st-order polynomial had two bases: x^1 and x^0 .

We could add another column, with the values of x^2 .
This will lead to a 2nd-order polynomial.

If we first make a matrix that looks like this

$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ \vdots & \vdots \end{bmatrix}$$

then multiply it by a vector $\begin{pmatrix} m \\ c \end{pmatrix}$

$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \end{bmatrix} \begin{pmatrix} m \\ c \end{pmatrix} =$$

ACTIVITY: What does this equal?

From slide 16...

Basis Functions

We've already used a basis, the 1st-order polynomial had two bases: x^1 and x^0 .

We could add another column, with the values of x^2 .
This will lead to a 2nd-order polynomial.

The solution is the same expression as before.

Our parameter vector would just need another number (to describe the contribution of the quadratic term).

If we first make a matrix that looks like this

$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \\ \vdots & \vdots \end{bmatrix}$$

then multiply it by a vector (c^m)

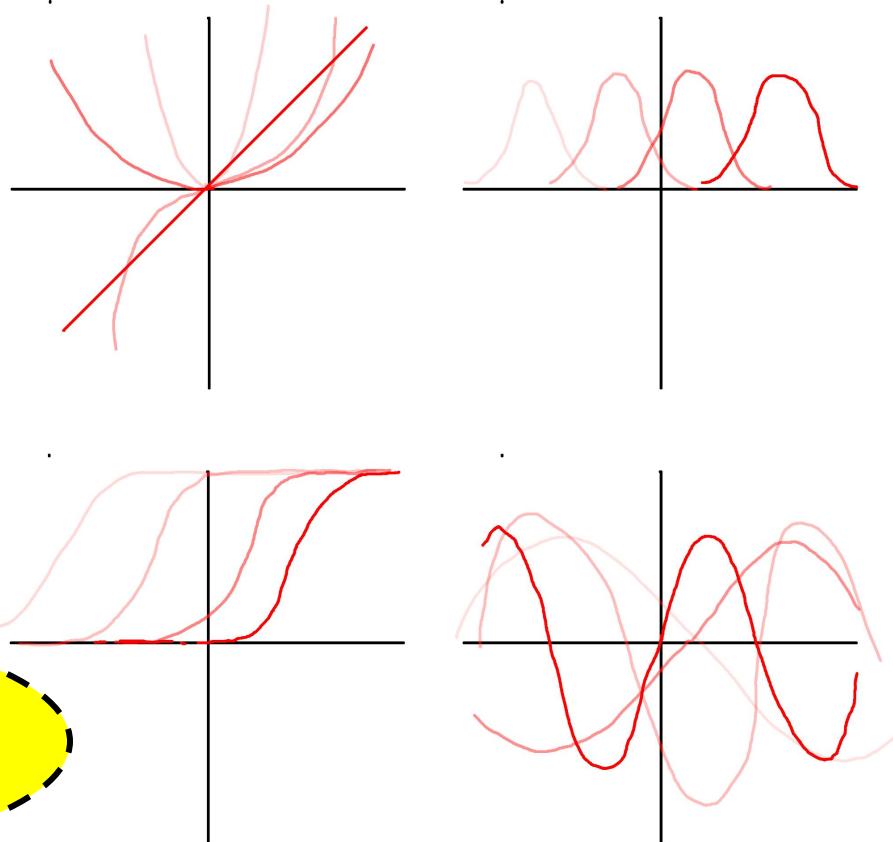
$$\begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ x_3 & 1 \end{bmatrix} \begin{bmatrix} c^m \end{bmatrix} =$$

ACTIVITY: What does this equal?

From slide 16...

Example Basis Functions

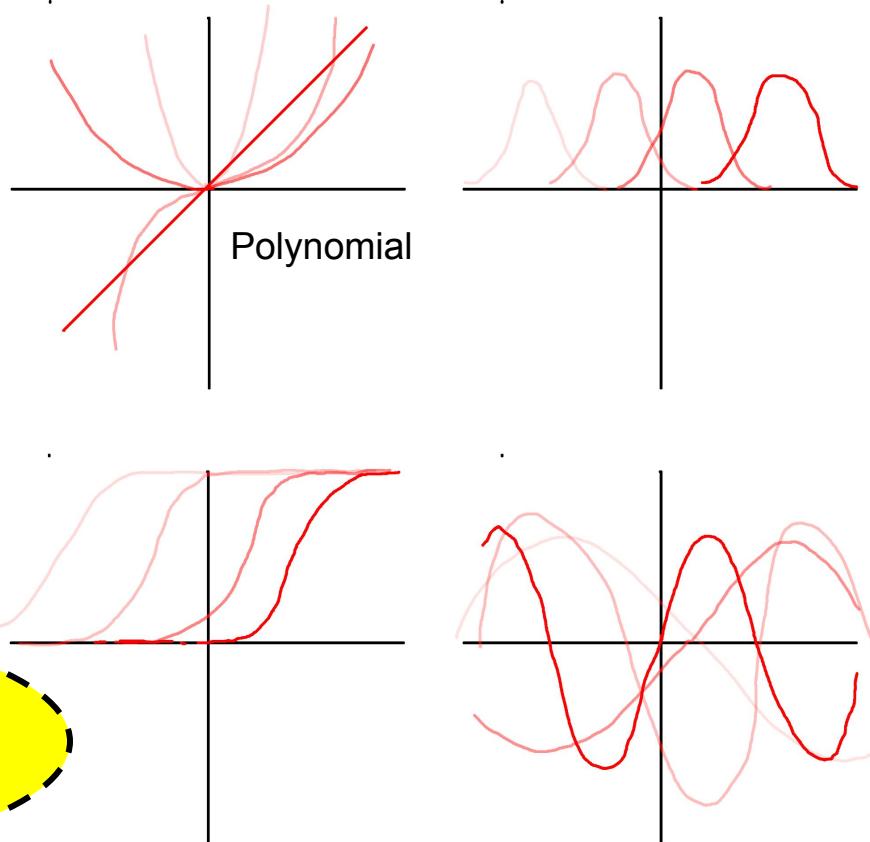
- Polynomial
- Gaussian
- Fourier
- Sigmoid



ACTIVITY: Match the words to the bases. One of them is a 'random' basis (so not picked from a grid)

Example Basis Functions

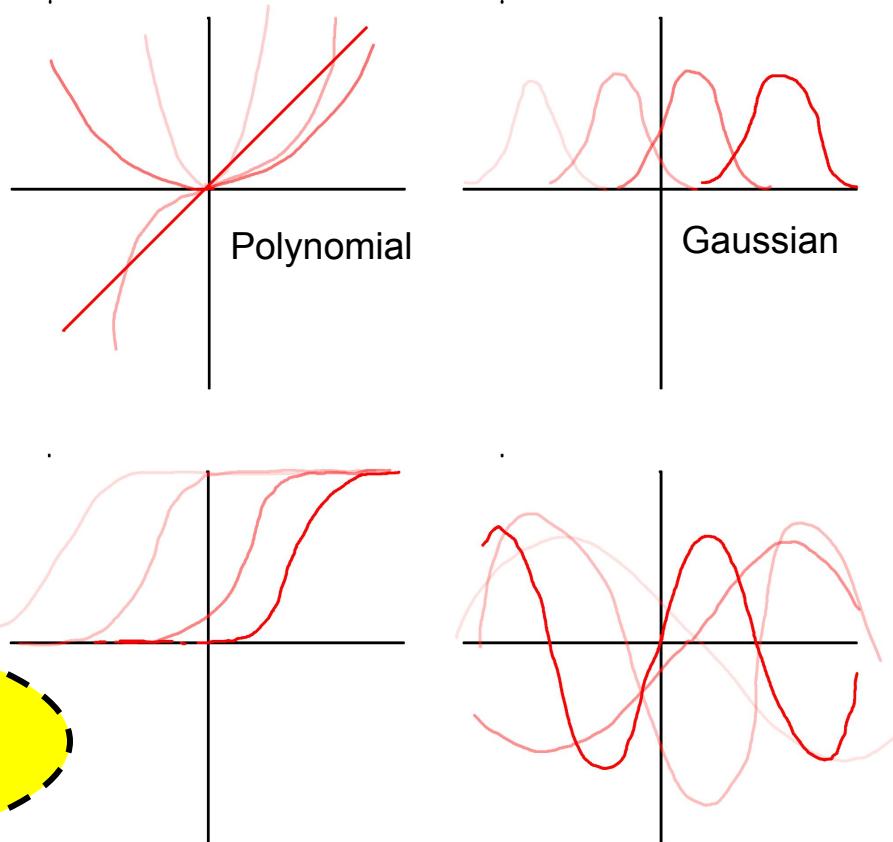
- Polynomial
- Gaussian
- Fourier
- Sigmoid



ACTIVITY: Match the words to the bases. One of them is a 'random' basis (so not picked from a grid)

Example Basis Functions

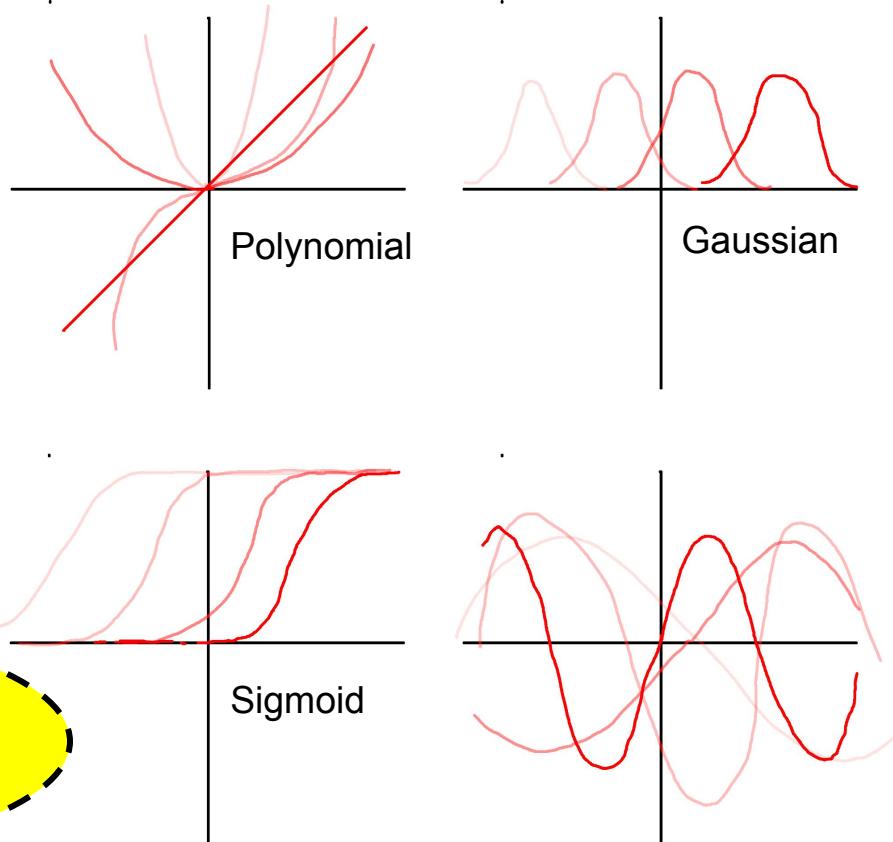
- Polynomial
- Gaussian
- Fourier
- Sigmoid



ACTIVITY: Match the words to the bases. One of them is a 'random' basis (so not picked from a grid)

Example Basis Functions

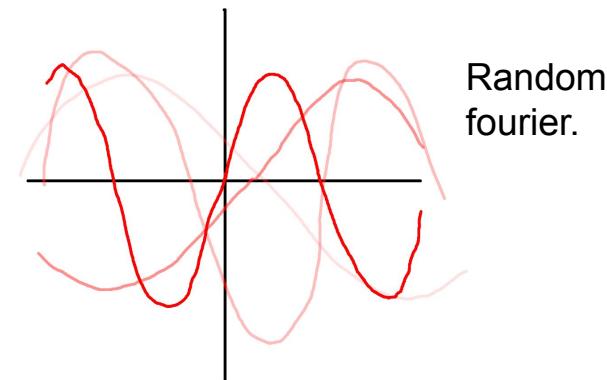
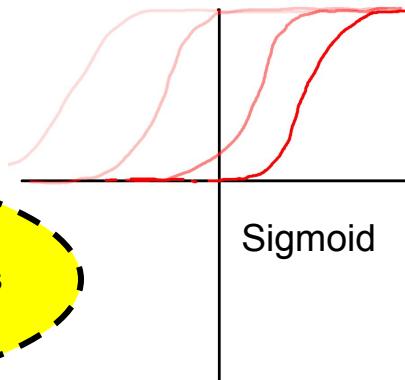
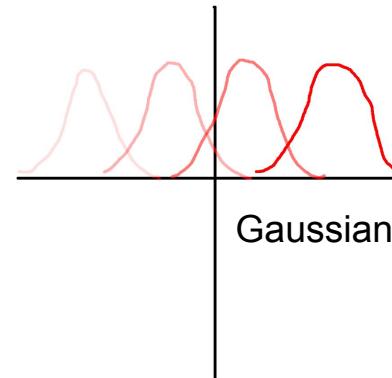
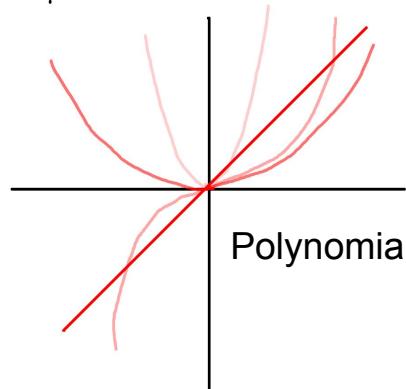
- Polynomial
- Gaussian
- Fourier
- Sigmoid



ACTIVITY: Match the words to the bases. One of them is a 'random' basis (so not picked from a grid)

Example Basis Functions

- Polynomial
- Gaussian
- Fourier
- Sigmoid



ACTIVITY: Match the words to the bases. One of them is a 'random' basis (so not picked from a grid)

Regularisation

A common problem is ‘overfitting’. You can get parameter values that are very large to explain data.

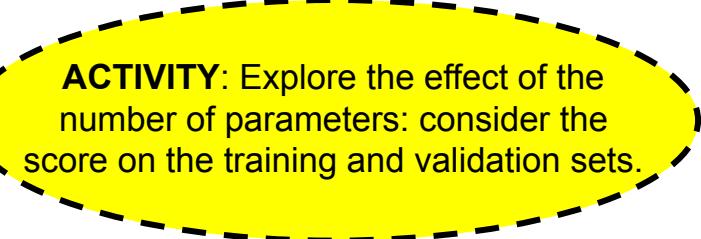
See [the example](#).

To help avoid this we can add an additional term to the cost function that penalises lots of large parameters.

$$h(w) = E(w) + \lambda R(w)$$

Annotations:

- E : Our original cost function
- λ : Regulariser parameter (maybe set using validation set?).
- R : regulariser



Regularisation

A common problem is ‘overfitting’. You can get parameter values that are very large to explain data.

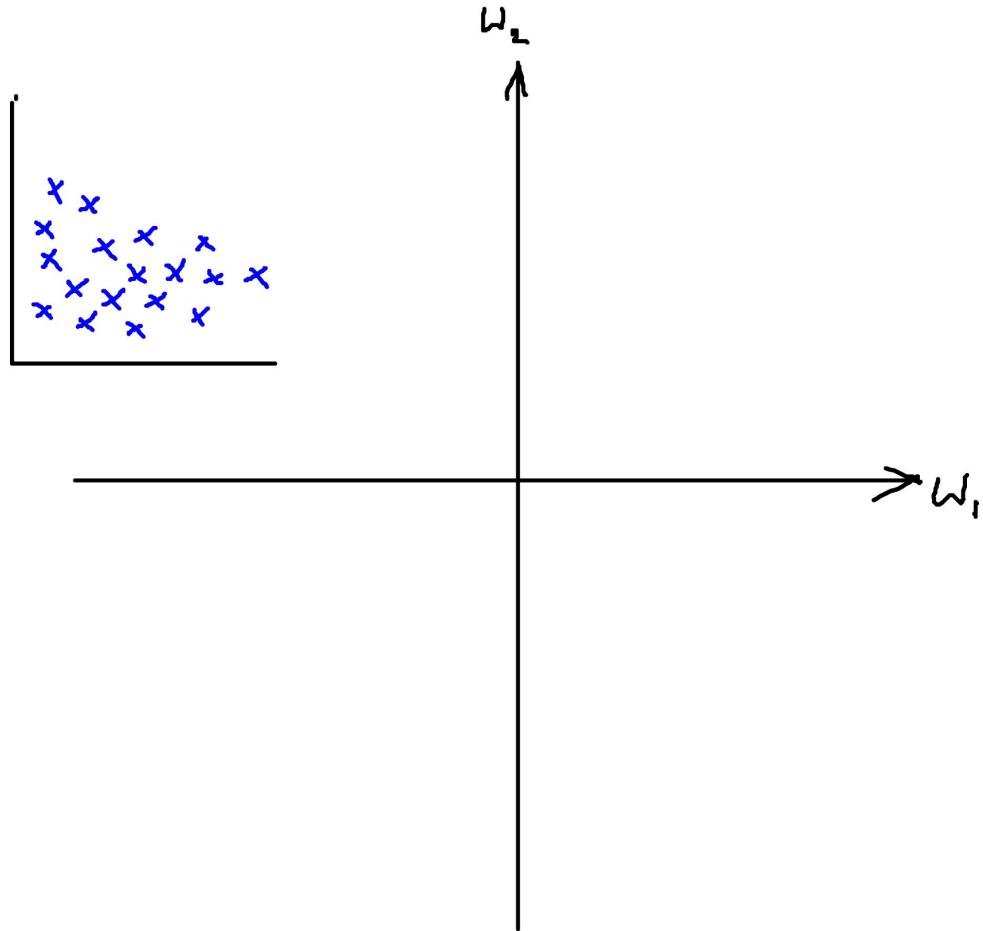
See [the example](#).

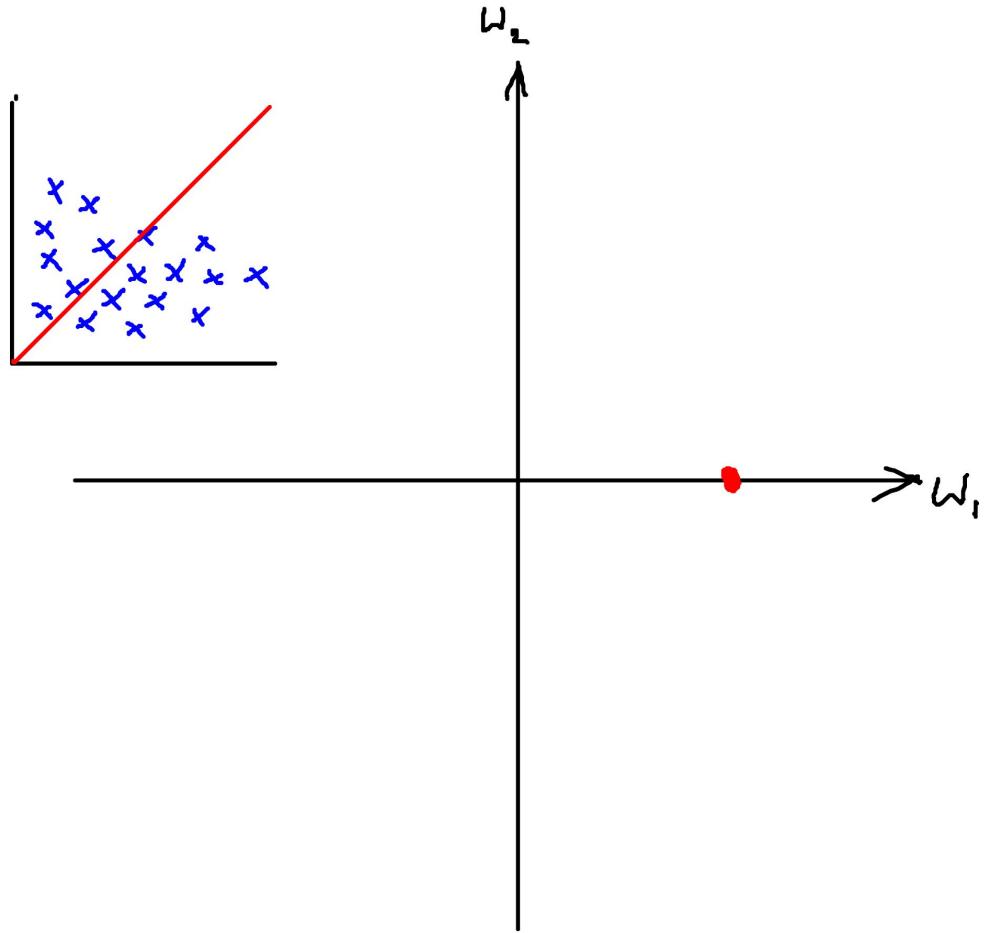
To help avoid this we can add an additional term to the cost function that penalises lots of large parameters.

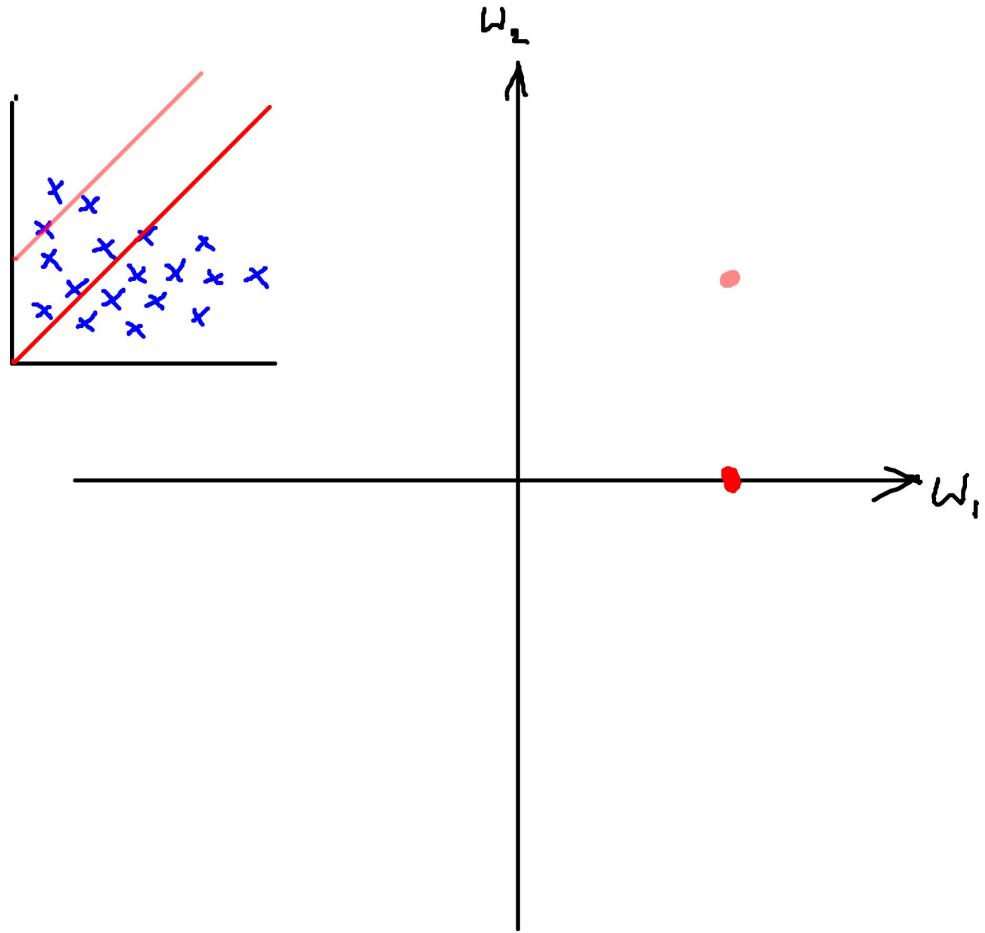
$$h(w) = E(w) + \lambda R(w)$$

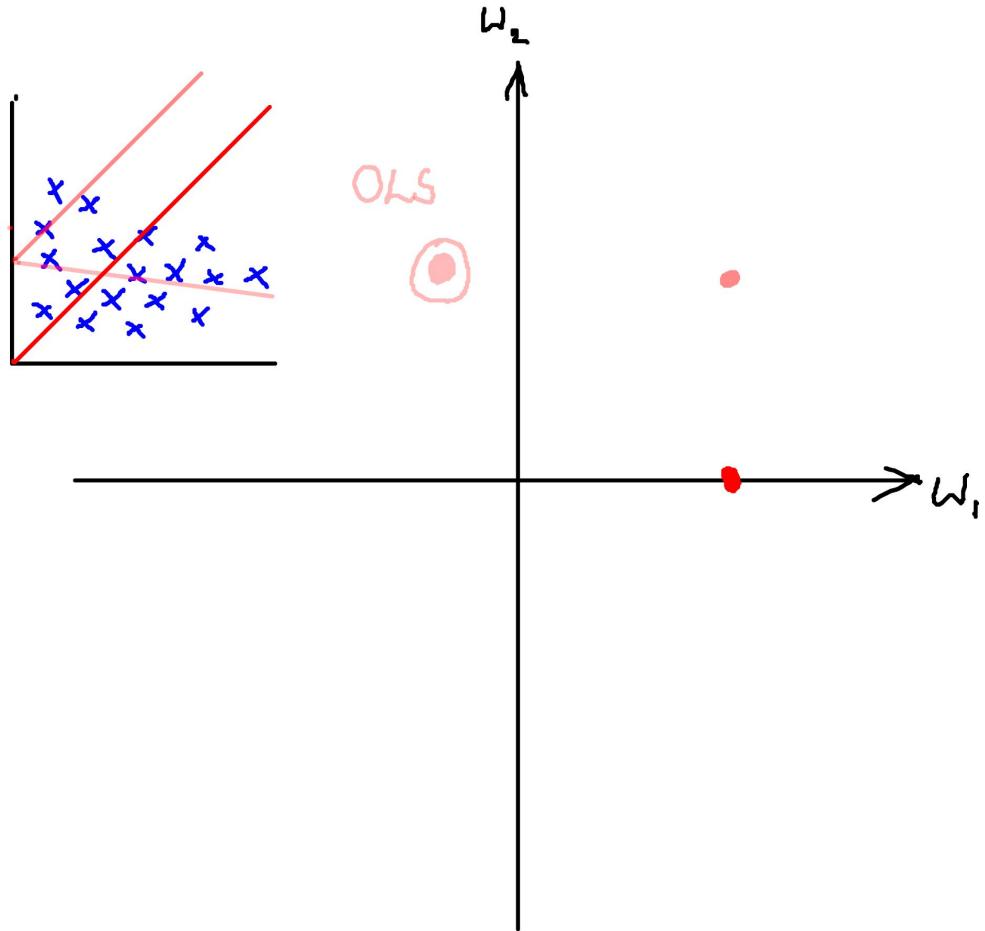
Annotations:

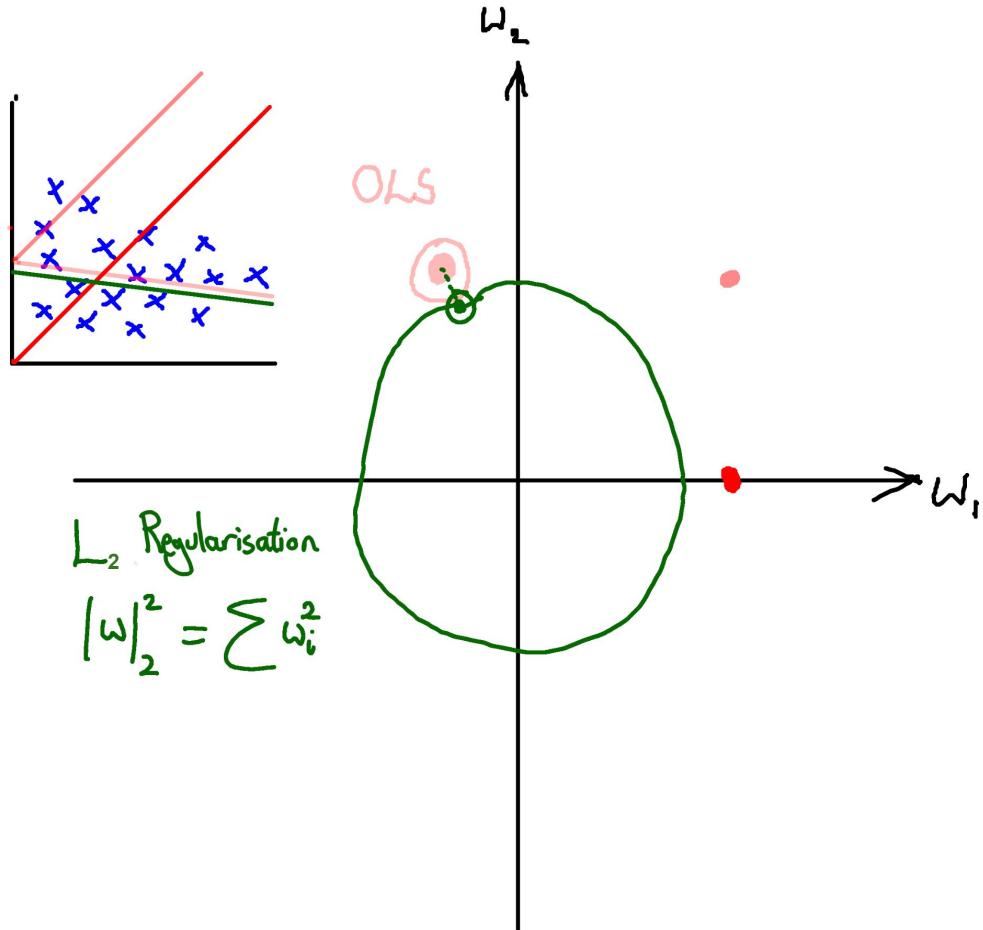
- E : Our original cost function
- λ : Regulariser parameter (maybe set using validation set?).
- R : regulariser





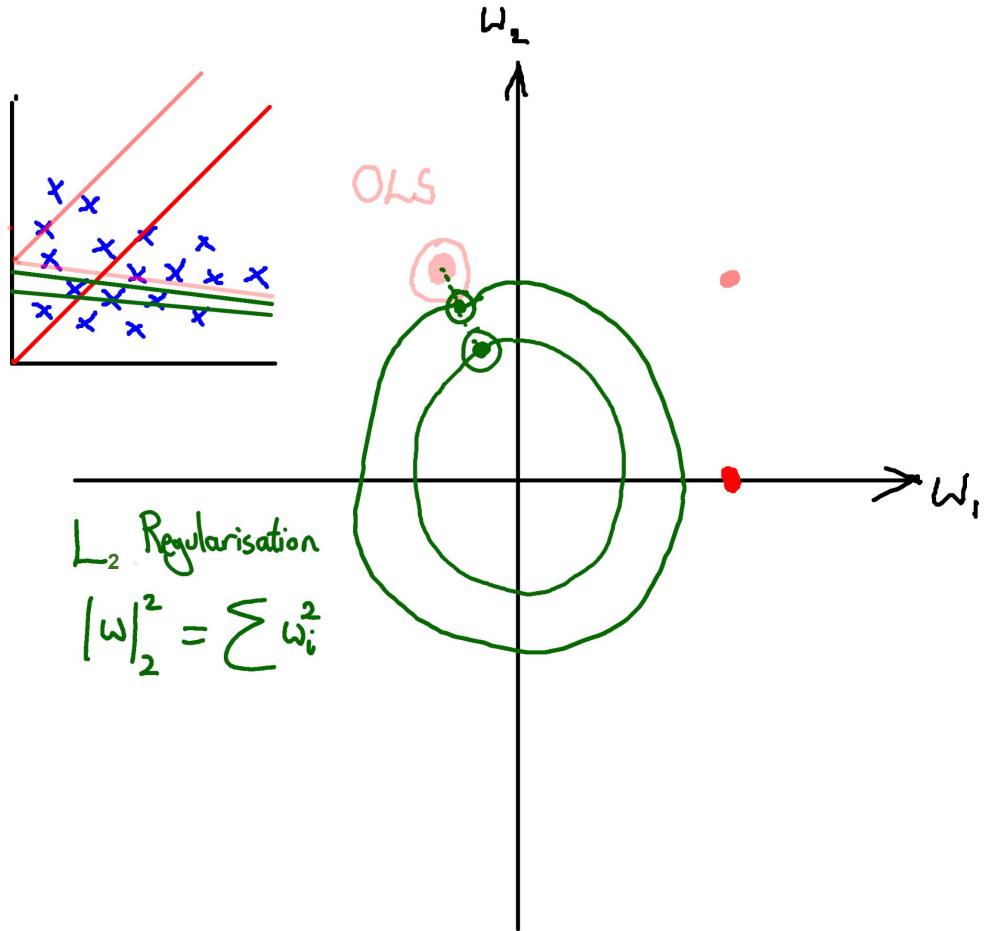


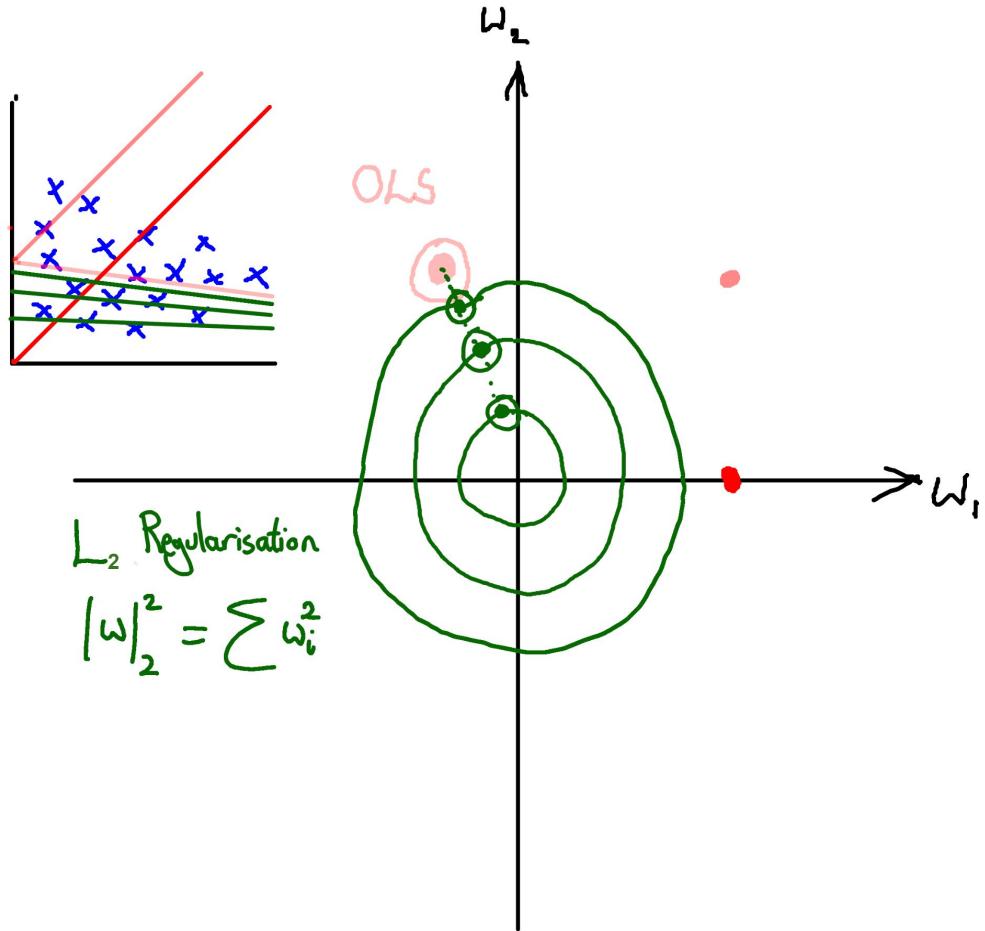


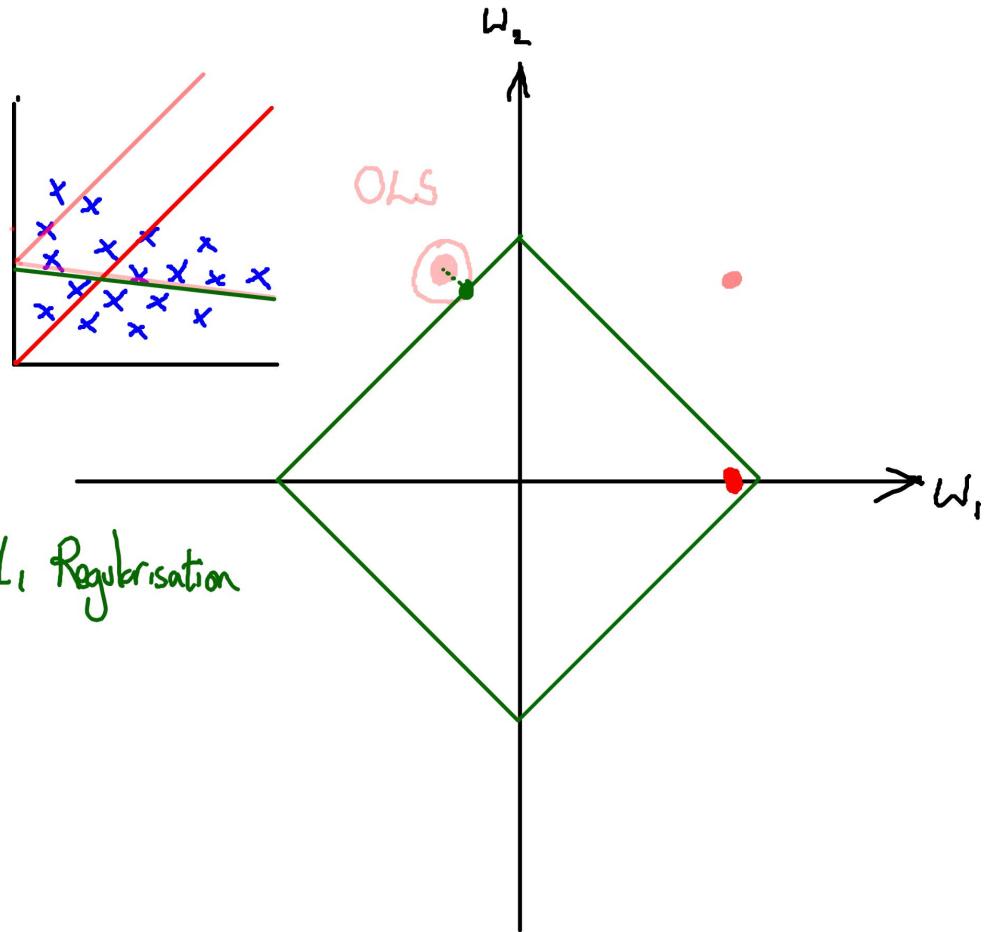


L_2 regularisation - also called ridge regression.

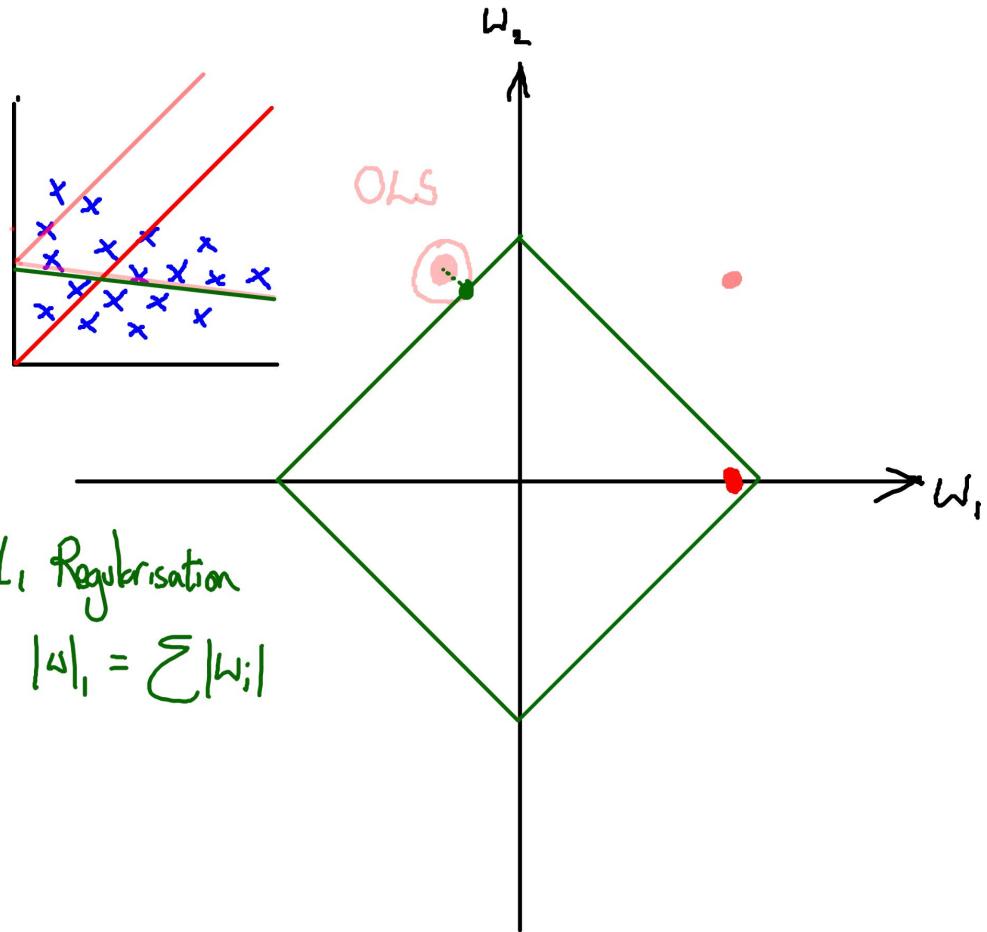
In this toy example we can't see much benefit. But with many bases (or parameters) this often is absolutely necessary to make the regression work.



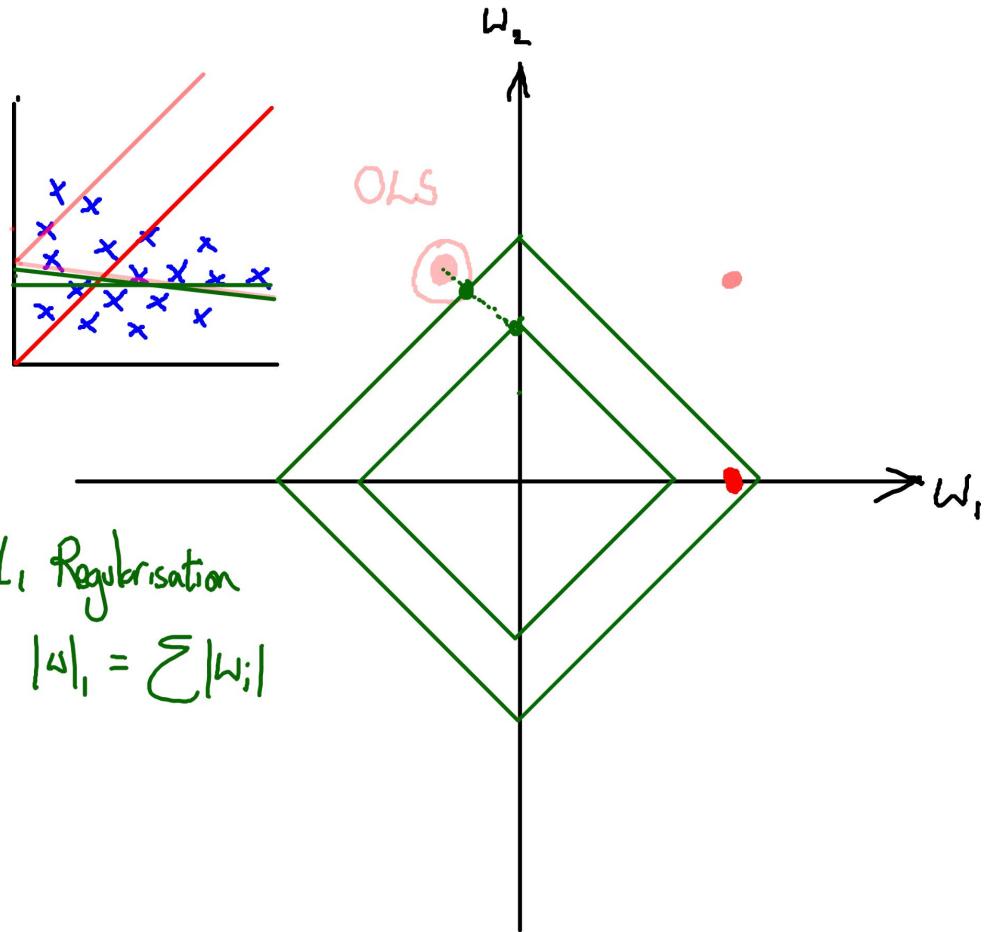




L_1 regularisation has the property that it will push many of the parameters to zero. This can be useful if you think only a few features are relevant.



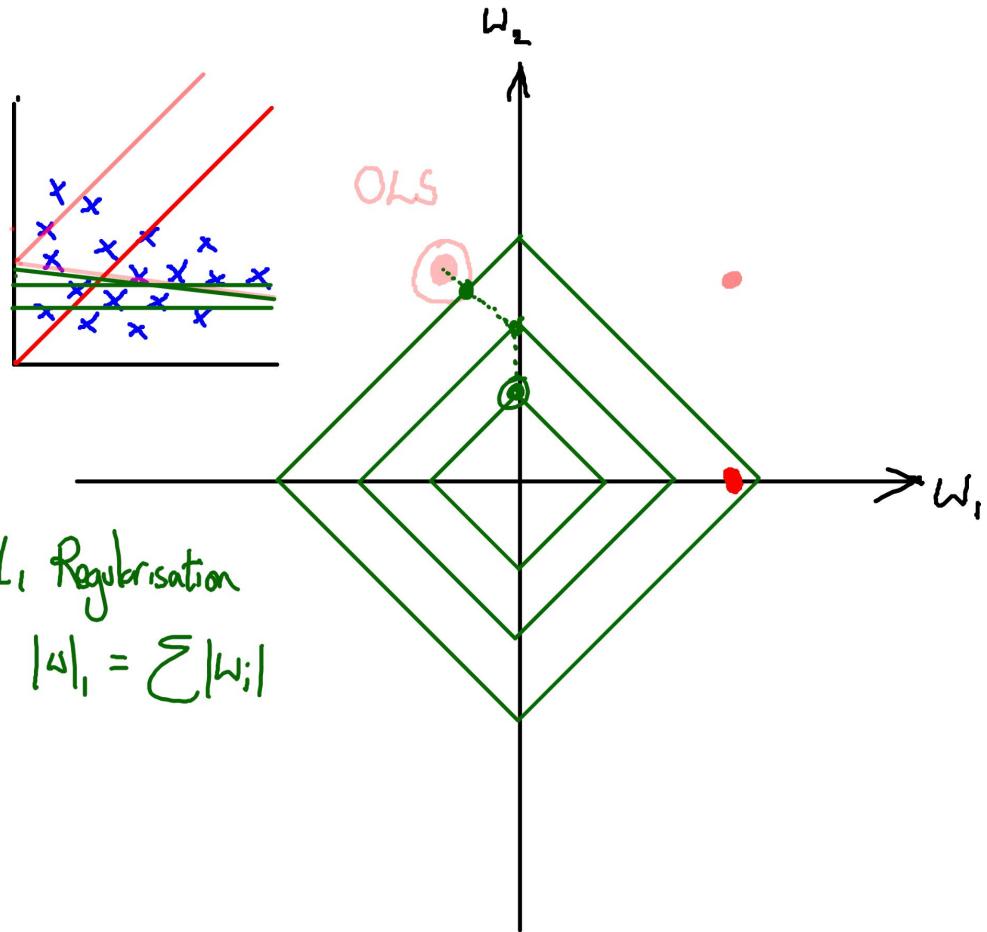
L_1 regularisation has the property that it will push many of the parameters to zero. This can be useful if you think only a few features are relevant.



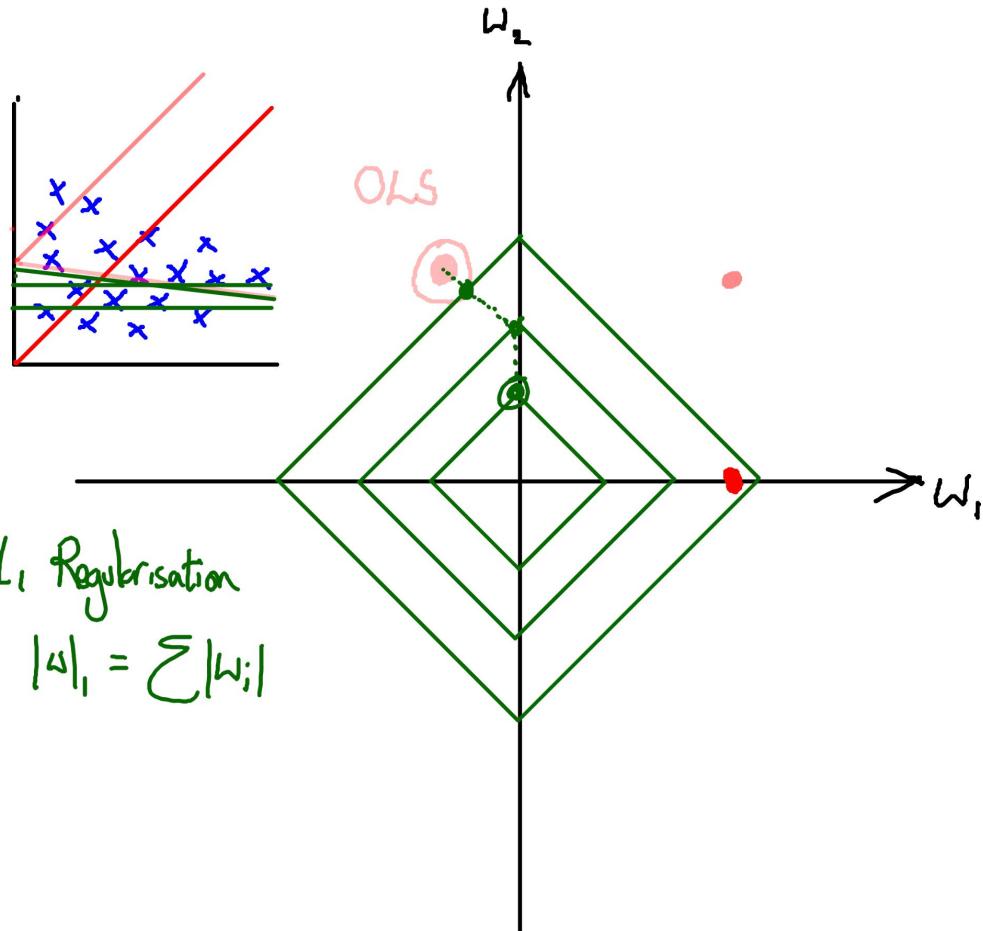
L_1 Regularisation

$$|w|_1 = \sum |w_i|$$

L_1 regularisation has the property that it will push many of the parameters to zero. This can be useful if you think only a few features are relevant.



L_1 regularisation has the property that it will push many of the parameters to zero. This can be useful if you think only a few features are relevant.



L_1 regularisation has the property that it will push many of the parameters to zero. This can be useful if you think only a few features are relevant.

Elastic net regularisation has some L_1 and some L_2 , mixed together.

Other approaches exist.

Summary / Take home messages

- We can write the sum square cost function as

$$E = (y - X\omega)^T (y - X\omega)$$

Summary / Take home messages

- We can write the sum square cost function as $E = (y - X\omega)^T (y - X\omega)$
- We can differentiate, set to zero, and find a closed form solution for linear regression.

Summary / Take home messages

- We can write the sum square cost function as $E = (y - X\omega)^T (y - X\omega)$
- We can differentiate, set to zero, and find a closed form solution for linear regression.
- This is the same solution as when we use the maximum likelihood approach.

Summary / Take home messages

- We can write the sum square cost function as $E = (y - X\omega)^T (y - X\omega)$
- We can differentiate, set to zero, and find a closed form solution for linear regression.
- This is the same solution as when we use the maximum likelihood approach.
- If we don't have a closed form solution, we need to optimise our parameters iteratively.

Summary / Take home messages

- We can write the sum square cost function as $E = (y - X\omega)^T (y - X\omega)$
- We can differentiate, set to zero, and find a closed form solution for linear regression.
- This is the same solution as when we use the maximum likelihood approach.
- If we don't have a closed form solution, we need to optimise our parameters iteratively.
- Steepest descent is simple but often is slow.

Summary / Take home messages

- We can write the sum square cost function as $E = (y - X\omega)^T (y - X\omega)$
- We can differentiate, set to zero, and find a closed form solution for linear regression.
- This is the same solution as when we use the maximum likelihood approach.
- If we don't have a closed form solution, we need to optimise our parameters iteratively.
- Steepest descent is simple but often is slow.
- Other options: line search, Newton's method, conjugate gradient descent.

Summary / Take home messages

- We can write the sum square cost function as $E = (y - X\omega)^T (y - X\omega)$
- We can differentiate, set to zero, and find a closed form solution for linear regression.
- This is the same solution as when we use the maximum likelihood approach.
- If we don't have a closed form solution, we need to optimise our parameters iteratively.
- Steepest descent is simple but often is slow.
- Other options: line search, Newton's method, conjugate gradient descent.
- If we have too much data we could use mini-batching and stochastic gradient descent.

Summary / Take home messages

- We can write the sum square cost function as $E = (y - X\omega)^T (y - X\omega)$
- We can differentiate, set to zero, and find a closed form solution for linear regression.
- This is the same solution as when we use the maximum likelihood approach.
- If we don't have a closed form solution, we need to optimise our parameters iteratively.
- Steepest descent is simple but often is slow.
- Other options: line search, Newton's method, conjugate gradient descent.
- If we have too much data we could use mini-batching and stochastic gradient descent.
- If the model overfits, try regularising.