

# Lecture 3. Linear Regression

COMP90051 Statistical Machine Learning

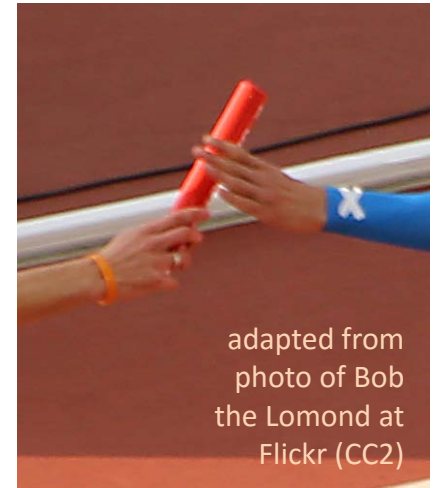
Semester 2, 2017  
Lecturer: Andrey Kan



THE UNIVERSITY OF  
MELBOURNE

# Weeks 2 to 8 inclusive

- Lecturer: Andrey Kan
  - \* MS: Moscow, PhD: Melbourne
  - \* Email: [andrey.kan@unimelb.edu.au](mailto:andrey.kan@unimelb.edu.au)
  - \* Office hour: after Thursday's lecture
- Past: software engineer
  - \* Optimising compilers
  - \* Multimedia framework
- Present: postdoc
  - \* Medical image analysis
  - \* Computational biology
- Future: applied scientist
  - \* Machine Learning
  - \* *Going to Seattle at the end of September!*



# This lecture

- Linear regression
  - \* Worked example and the model
  - \* Regression as a probabilistic model
- Regularisation
  - \* Irrelevant features and an ill-posed problem
  - \* Regulariser as a prior

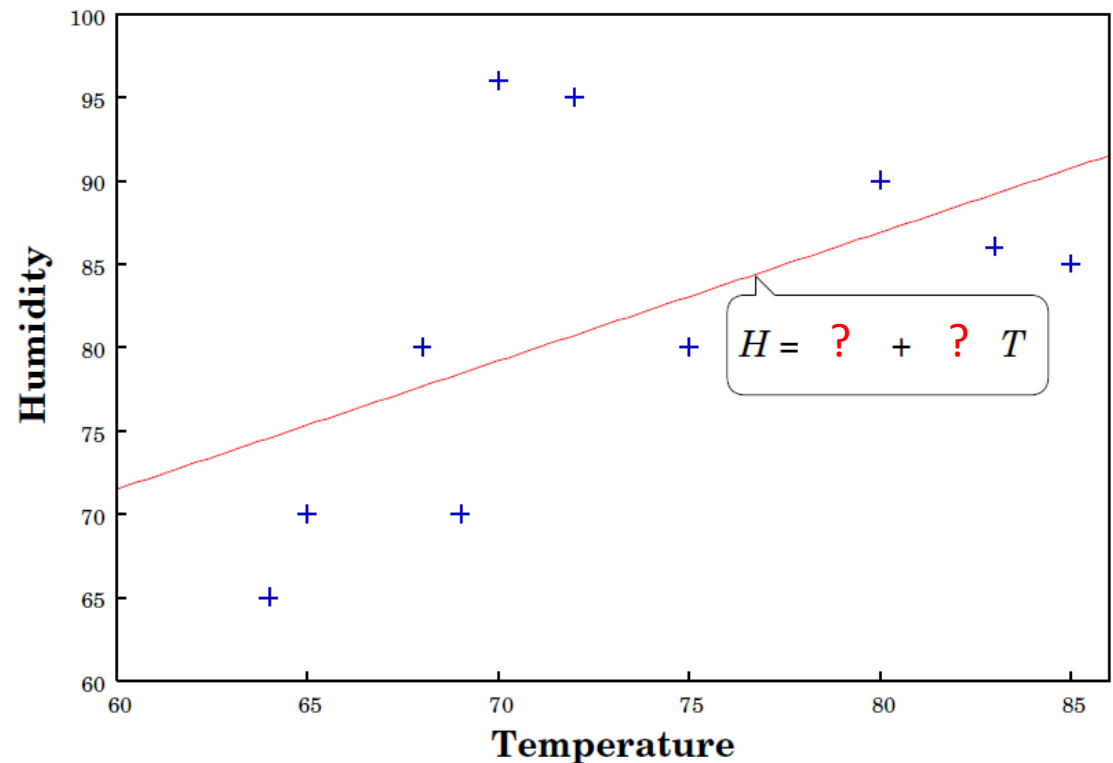
# Linear Regression Model

A simple model that tends to require less data and be easier to interpret.

It offers mathematical convenience at the expense of flexibility.

# Example: Predict humidity from temperature

Temperature	Humidity
TRAINING DATA	
85	85
80	90
83	86
70	96
68	80
65	70
64	65
72	95
69	70
75	80
TEST DATA	
75	70



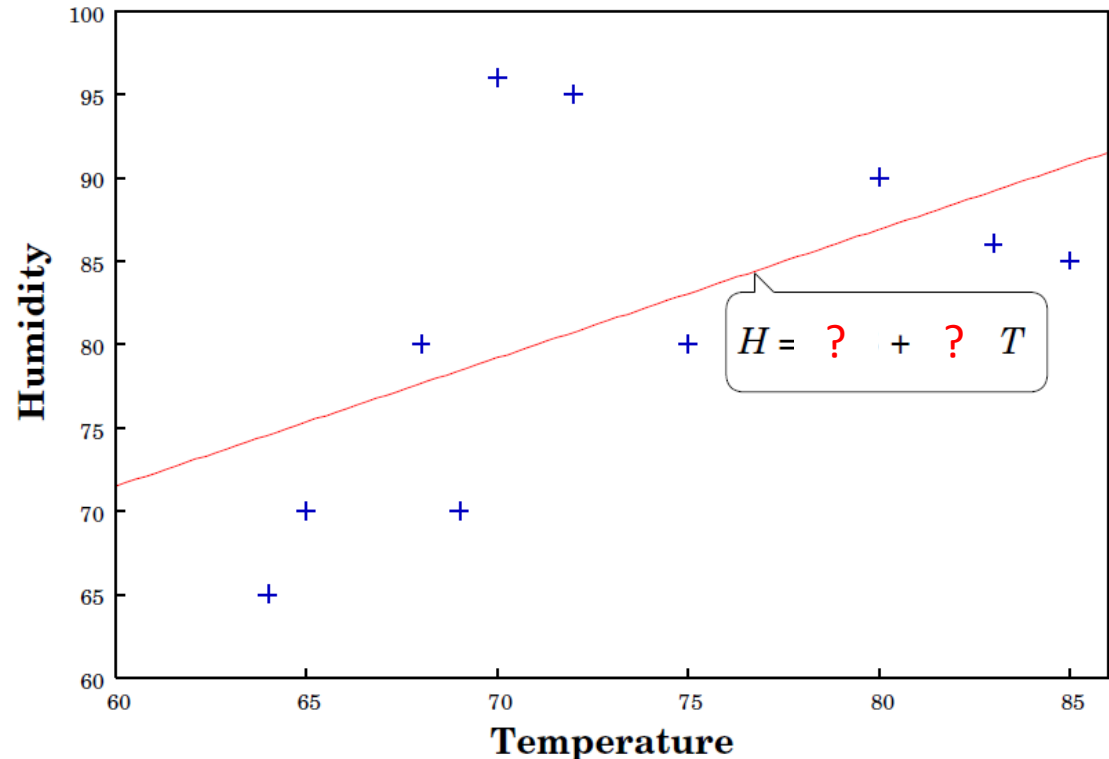
In regression, the task is to predict numeric response (*aka* dependent variable) from features (*aka* predictors or independent variables)

**Assume a linear relation:**  $H = a + bT$

( $H$  – humidity;  $T$  – temperature;  $a, b$  – parameters)

# Example: Problem statement

- The model is  
$$H = a + bT$$
- Fitting the model  
= finding “best”  
 $a, b$  values for  
data at hand
- Popular criterion:  
minimise the **sum  
of squared errors**  
(aka residual sum  
of squares)



# Example: Finding parameter values

To find  $a, b$  that minimise  $L = \sum_{i=1}^{10} (H_i - (a + b T_i))^2$

set derivatives to zero:

$$\frac{\partial L}{\partial a} = -2 \sum_{i=1}^{10} (H_i - a - b T_i) = 0$$

if we know  $b$ , then  $\hat{a} = \frac{1}{10} \sum_{i=1}^{10} (H_i - b T_i)$

$$\frac{\partial L}{\partial b} = -2 \sum_{i=1}^{10} T_i (H_i - a - b T_i) = 0$$

if we know  $a$ , then  $\hat{b} = \frac{1}{\sum_{i=1}^{10} T_i^2} \sum_{i=1}^{10} T_i (H_i - a)$

Basic calculus:

- Write derivative
- Set to zero
- Solve for model

**Coordinate descent:** guess  $a$ , solve for  $b$ ; solve for  $a$ ; repeat gives linear regression  $a = 25.3, b = 0.77$  (requires many iterations!)

# Example: Analytic solution

- Can we do better? We have two equations and two unknowns  $a, b$
- Rewrite as a system of linear equations

$$\begin{pmatrix} 10 & \sum_{i=1}^{10} T_i \\ \sum_{i=1}^{10} T_i & \sum_{i=1}^{10} T_i^2 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^{10} H_i \\ \sum_{i=1}^{10} T_i H_i \end{pmatrix}$$

- **Analytic solution:**  $a = 25.3, b = 0.77$
- (Can solve using `numpy.linalg.solve` or equivalent)



# Linear regression model

- Assume a linear relationship between response  $y \in \mathbb{R}$  and an instance with features  $x_1, \dots, x_m \in \mathbb{R}$

$$y \approx w_0 + \sum_{i=1}^m x_i w_i$$

Here  $w_1, \dots, w_m \in \mathbb{R}$  denote weights (model parameters)

- **Trick:** add a dummy feature  $x_0 = 1$  and use vector notation

$$y \approx \sum_{i=0}^m x_i w_i = \mathbf{x}' \mathbf{w}$$

# Checkpoint

- Consider a dataset  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  and some parameter vector  $\mathbf{w}$  (same dimensionality as  $\mathbf{x}_i$ ). Which of the following statements is necessarily true



Each  $y_i$  can be expressed as  $(\mathbf{x}_i)' \mathbf{w}$

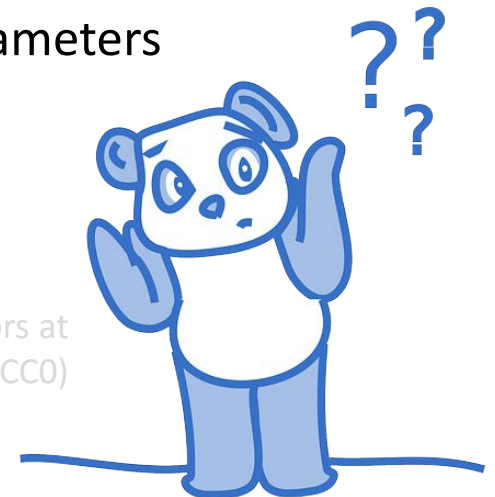


Given  $\mathbf{w}$ , it is always possible to compute the sum of squared errors for this dataset



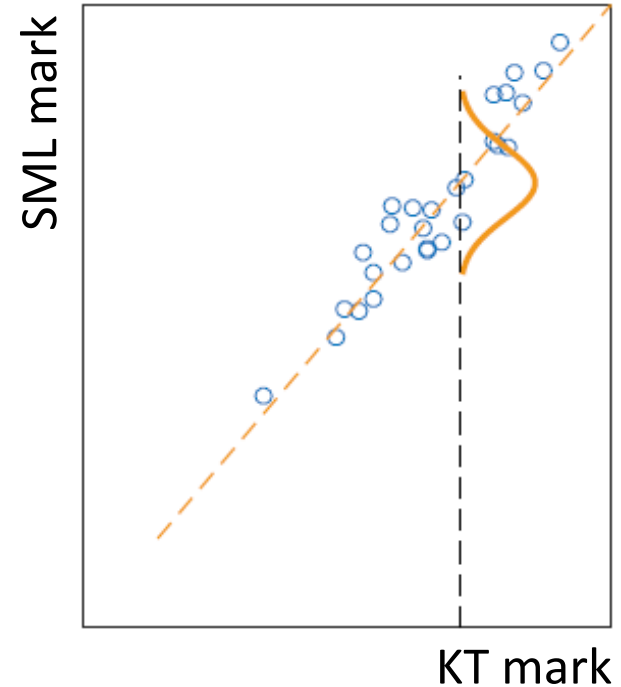
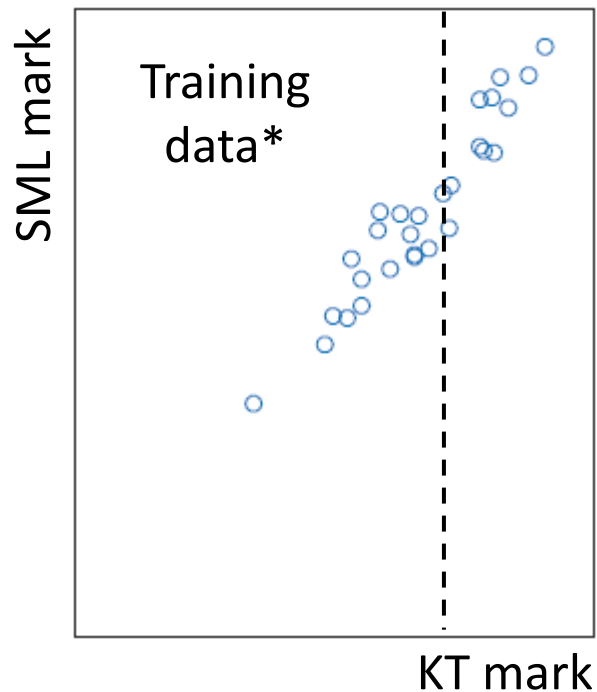
Linear regression model for this data has  $n$  parameters

art: OpenClipartVectors at  
pixabay.com (CC0)



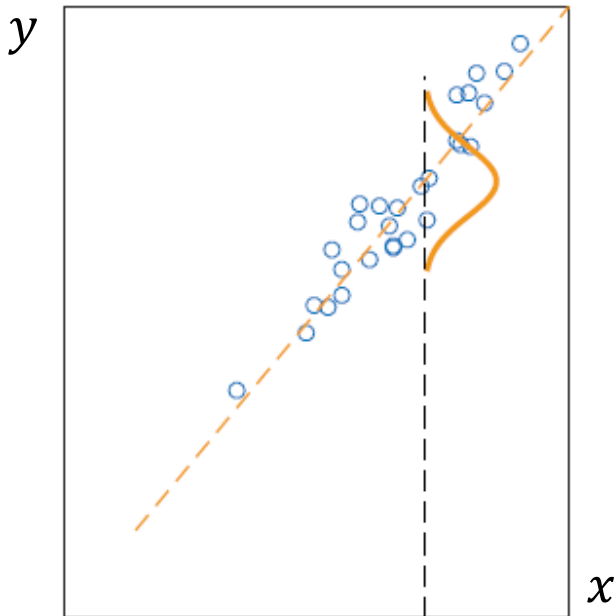
# Data is noisy!

Example: predict mark for Statistical Machine Learning (SML)  
from mark for Knowledge Technologies (KT)



\* synthetic data :)

# Regression as a probabilistic model



- Assume a **probabilistic model**:  
 $\mathcal{Y} = \mathbf{x}'\mathbf{w} + \varepsilon$ 
  - \* Here  $\mathcal{Y}$  and  $\varepsilon$  are random variables
  - \* Variable  $\varepsilon$  encodes noise
- Next, assume normally distributed noise:  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

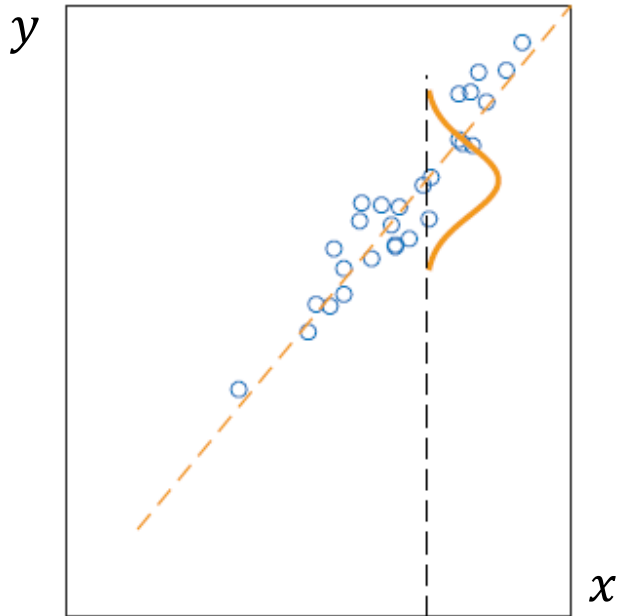
- Recall that  $\mathcal{N}(x|\mu, \sigma^2) \equiv \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

- Therefore

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

this is a squared error!

# Parametric probabilistic model



- Using simplified notation, our model is

$$p(y|\mathbf{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mathbf{x}'\mathbf{w})^2}{2\sigma^2}\right)$$

- Note that parameters are now  $\mathbf{w}, \sigma^2$

- Given observed data  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ , we want to find parameter values that “best” explain the data
- **Maximum likelihood estimation**: choose parameter values that maximise the probability of observed data

# Maximum likelihood estimation

- Assuming independence of data points, the probability of data is

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n p(y_i | \mathbf{x}_i)$$

- Recall that  $p(y_i | \mathbf{x}_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - (\mathbf{x}_i)' \mathbf{w})^2}{2\sigma^2}\right)$

- “Log trick”: Instead of maximising this quantity, maximise the logarithm

$$\sum_{i=1}^n \log p(y_i | \mathbf{x}_i) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\mathbf{x}_i)' \mathbf{w})^2 + K$$

here  $K$  doesn't depend on  $\mathbf{w}$

the sum of squared errors!

- Under this model, maximising log-likelihood as a function of  $\mathbf{w}$  is equivalent to minimising the sum of squared errors

# Method of least squares

- Training data:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ . Note bold face in  $\mathbf{x}_i$
- For convenience, place instances in rows (so attributes go in columns), representing training data as an  $n \times (m + 1)$  matrix  $\mathbf{X}$ , and  $n \times 1$  vector  $\mathbf{y}$
- **The model** assumes  $\mathbf{y} \approx \mathbf{X}\mathbf{w}$
- To find  $\mathbf{w}$ , minimise the **sum of squared errors**

$$L = \sum_{i=1}^n \left( y_i - \sum_{j=0}^m X_{ij} w_j \right)^2$$

Basic calculus:

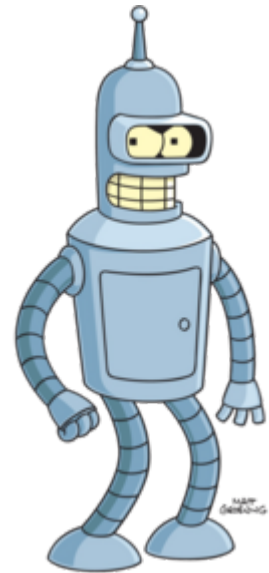
- Write derivative
- Set to zero
- Solve for model

- **Setting derivative to zero** and solving for  $\mathbf{w}$  yields
 
$$\hat{\mathbf{w}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

- \* This is a system of equations called the normal equations
- \* This system is defined only if the inverse exists

# Heads up!

- This subject covers a large variety of computational methods
- But there will be several recurring topics, common threads run throughout the entire course
- These topics reflect fundamental aspects of machine learning
- Basis expansion, representation
- Optimisation, loss functions
- Regularisation, overfitting





# Regularisation

Process of introducing additional information in order to solve an ill-posed problem or to prevent overfitting (*Wikipedia*)

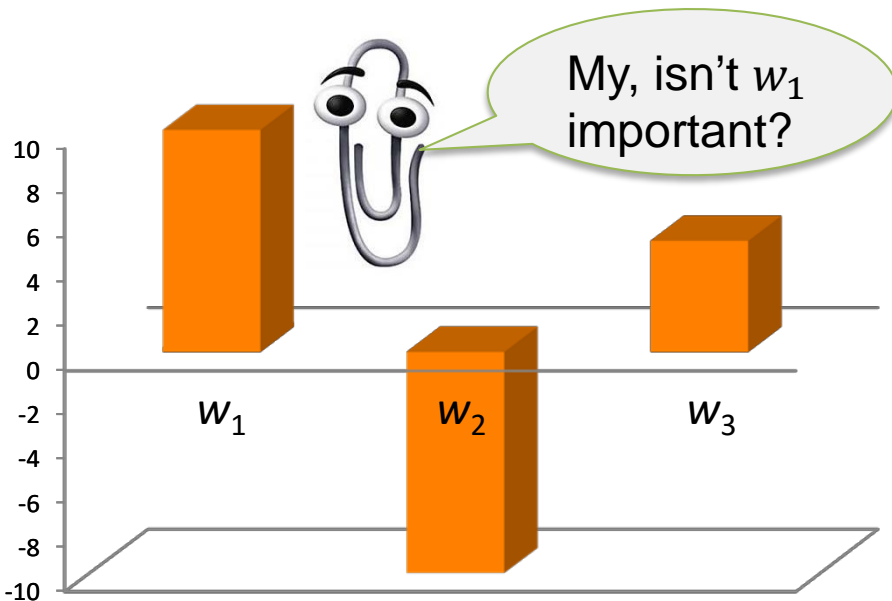
# Regularisation

- Major technique, common in Machine Learning
- Addresses one or more of the following related problems
  - \* Avoid ill-conditioning
  - \* Introduce prior knowledge
  - \* Constrain modelling
- This is achieved by augmenting the objective function
- In this lecture: we cover the first two aspects. We will cover more of regularisation throughout the subject

# Example 1: Irrelevant features

- Linear model on three features, first two same
  - \*  $\mathbf{X}$  is matrix a for  $n = 4$  instances (rows)
  - \* First two columns of  $\mathbf{X}$  identical
  - \* Feature 2 is **irrelevant** (or feature 1)
  - \* Model:  $y = w_1x_1 + w_2x_2 + w_3x_3 + w_0$

3	3	7
6	6	9
21	21	79
34	34	2



- Effect of perturbations on model predictions?
  - \* Add  $\Delta$  to  $w_1$
  - \* Subtract  $\Delta$  from  $w_2$
  - ...identical predictions
  - ...no interpretability

# Problems with irrelevant features

- In our example, suppose  $[\hat{w}_0, \hat{w}_1, \hat{w}_2, \hat{w}_3]$  is the “best solution”
- But for arbitrary  $\delta$  solution  $[\hat{w}_0, \hat{w}_1 + \delta, \hat{w}_2 - \delta, \hat{w}_3]'$  will lead to the same predictions and to the same sum of squared errors
- The solution is not unique
- One problem is lack of interpretability
- A more serious problem is that the finding the best parameters becomes an **ill-posed problem**

# Irrelevant (co-linear) features in general

- X-treme case: features complete clones
- For linear models, more generally
  - \* Feature  $X_{.j}$  is irrelevant if
  - \*  $X_{.j}$  is a **linear combination** of other columns

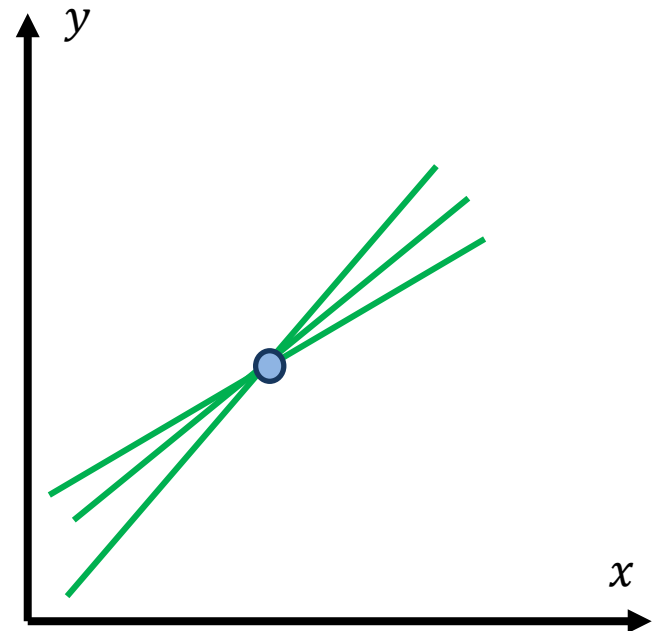
$$X_{.j} = \sum_{l \neq j} \alpha_l X_{.l}$$

... for some constants  $\alpha_l$

- Even *near*-irrelevance can be problematic
- Not just a pathological x-treme; ***easy to happen!***

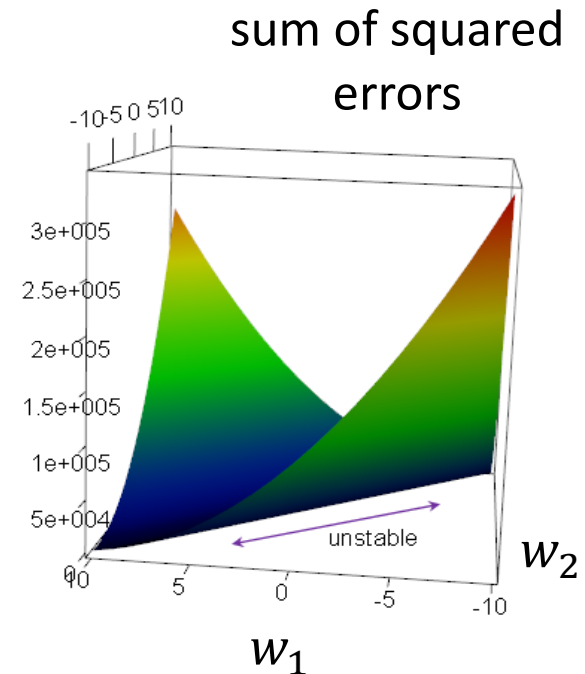
## Example 2: Lack of data

- Model is more complex than data
- Extreme example:
  - \* Model has two parameters (slope and intercept)
  - \* Only one data point
- Underdetermined system



# Ill-posed problems

- In both examples, finding the best parameters becomes an **ill-posed problem**
- This means that the problem solution is not defined
  - \* In our case  $w_1$  and  $w_2$  cannot be uniquely identified
- Remember the normal equations solution  $\hat{\mathbf{w}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$
- With irrelevant features,  $\mathbf{X}'\mathbf{X}$  has no inverse
- The system of linear equations has more unknowns than equations



convex, but not strictly convex

# Side note: L1 and L2 norms

- Throughout the course we will often encounter **norms** that are functions  $\mathbb{R}^n \rightarrow \mathbb{R}$  of a particular form
  - \* Intuitively, norms measure lengths of vectors in some sense

- More specifically, we will often use the L2 norm (*aka* **Euclidean distance**)

$$\|\mathbf{a}\| = \|\mathbf{a}\|_2 \equiv \sqrt{a_1^2 + \dots + a_n^2}$$

- And also the L1 norm (*aka* absolute norm or **Manhattan distance**)

$$\|\mathbf{a}\|_1 \equiv |a_1| + \dots + |a_n|$$

- For example, the sum of squared errors is a squared norm

$$L = \sum_{i=1}^n \left( y_i - \sum_{j=0}^m X_{ij} w_j \right)^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$



# Re-conditioning the problem

- Regularisation: introduce an **additional condition** into the system
- The original problem is to minimise  $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$
- The regularised problem is to minimise

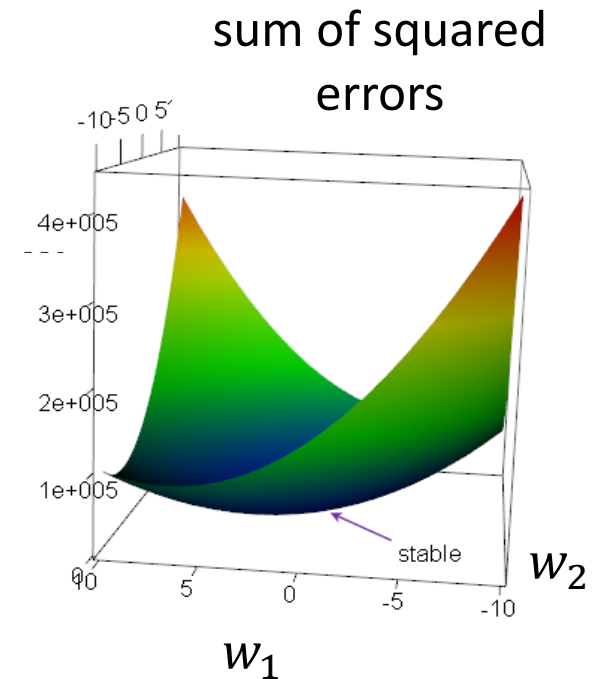
$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \text{ for } \lambda > 0$$

- The solution is now

$$\hat{\mathbf{w}} = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}'\mathbf{y}$$



- This formation is called **ridge regression**
  - \* Turns the ridge into a peak



# Regulariser as a prior

- Without regularisation model parameters are found based entirely on the information contained in the training set  $\mathbf{X}$
- Regularisation essentially means introducing additional information
- Recall our probabilistic model  $\mathcal{Y} = \mathbf{x}'\mathbf{w} + \varepsilon$ 
  - \* Here  $\mathcal{Y}$  and  $\varepsilon$  are random variables, where  $\varepsilon$  denotes noise
- Now suppose that  $\mathbf{w}$  is also a random variable (denoted as  $\mathcal{W}$ ) with a normal prior distribution

$$\mathcal{W} \sim \mathcal{N}(0, \lambda^2)$$

# Regulariser as a prior

- Now suppose that  $\mathbf{w}$  is also a random variable (denoted as  $\mathcal{W}$ ) with a normal prior distribution

$$\mathcal{W} \sim \mathcal{N}(0, \lambda^2)$$

- **Prior** = our initial expectations before seeing data
- In the above prior, we expect small weights and that no one feature dominates
  - \* Is this always appropriate? Consider data centring and scaling
- We could encode much more elaborate problem knowledge

# Computing posterior using Bayes rule

- The prior is then used to compute the posterior

Diagram illustrating the Bayes rule equation with callouts for each term:

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

Callouts:

- posterior (red callout pointing to  $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ )
- likelihood (green callout pointing to  $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$ )
- prior (purple callout pointing to  $p(\mathbf{w})$ )
- marginal likelihood (blue callout pointing to  $p(\mathbf{y}|\mathbf{X})$ )

- Instead of maximum likelihood (MLE), take **maximum a posteriori** estimate (MAP)
- Apply log trick, so that
 
$$\log(\text{posterior}) = \log(\text{likelihood}) + \log(\text{prior}) - \log(\text{marg})$$
- Arrive at the problem of minimising
 
$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda\|\mathbf{w}\|_2^2$$

this term doesn't affect optimisation

# This lecture

- Linear regression
  - \* Worked example and the model
  - \* Regression as a probabilistic model
- Regularisation
  - \* Irrelevant features and an ill-posed problem
  - \* Regulariser as a prior