

Machine Learning, Machine Learning (extended)

3 - Supervised Learning:
Linear Modelling by Maximum Likelihood
Kashif Rajpoot

k.m.rajpoot@cs.bham.ac.uk

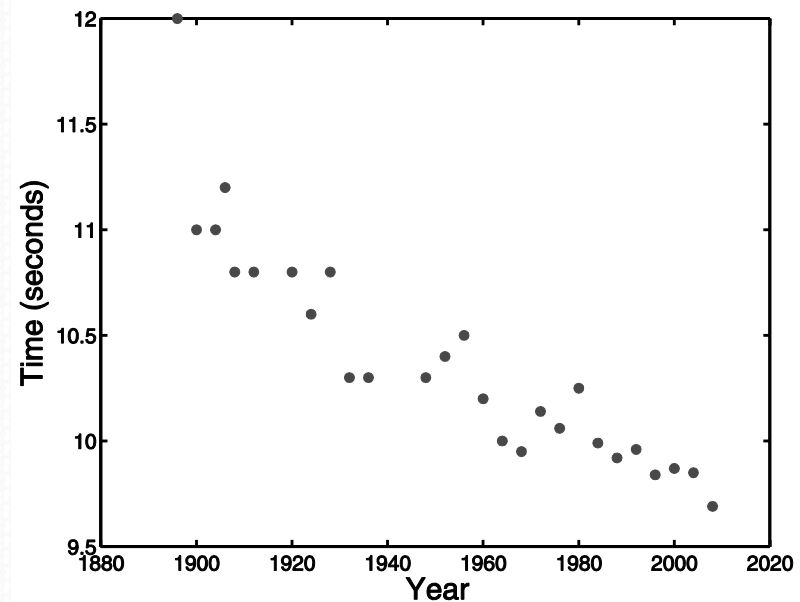
School of Computer Science
University of Birmingham

Outline

- Linear modelling
- Error = noise
- Thinking generatively
- Likelihood
- Maximum likelihood
- Model complexity
- Bias-variance tradeoff

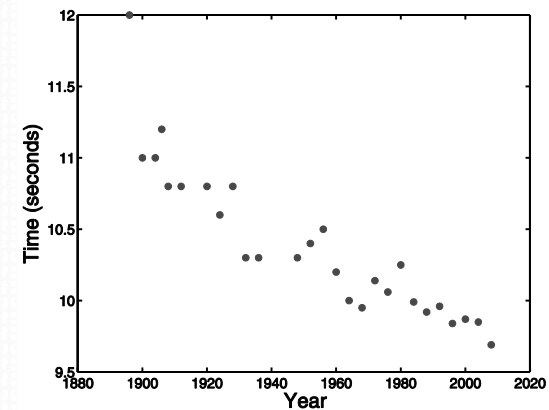
Linear modelling

- One of the most straightforward learning problems
 - Learn a linear function between attributes and responses
- Is there a functional dependence between Olympics year and 100m winning time?
 - Draw a line?
- Can we predict winning time for future games?



Linear modelling

- Learner model/function
 - Maps input attributes to output response
- Let's consider we can predict time $t = f(x)$
 - x ?
 - t ?
- Training samples
 - N attribute-response pairs
 $(x_1, t_1), (x_2, t_2), \dots (x_N, t_N)$



Linear modelling

- Linear modelling by minimizing loss

$$t_n = \hat{w}_0 + \hat{w}_1 x_n$$

where the model parameters are estimated from Olympics data

$$\hat{w}_1 = \frac{\bar{x}\bar{t} - \bar{\bar{x}\bar{t}}}{\bar{x^2} - (\bar{x})^2}$$

$$\hat{w}_0 = \bar{t} - \hat{w}_1 \bar{x}$$

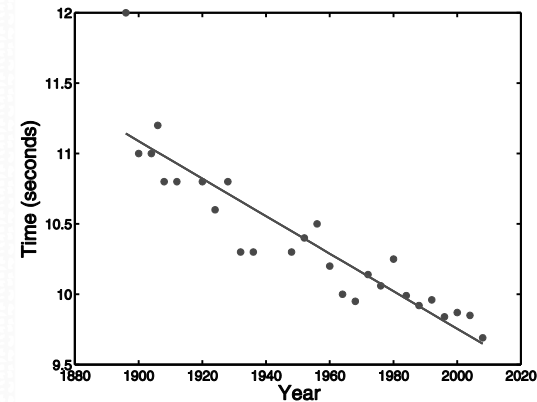
- With the vector notation

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

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \text{ and } \mathbf{x}_n = \begin{bmatrix} 1 \\ x_n \end{bmatrix}$$

while model parameters are estimated as: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$

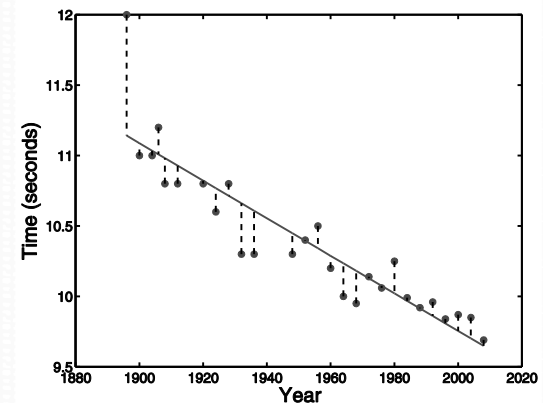
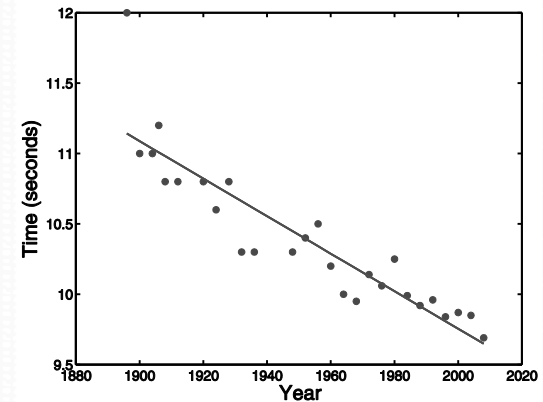
$$f(x; w_0, w_1) = 36.416 - 0.013x$$



$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

Linear modelling

- Linear modelling by minimizing loss
 - Model shows to capture trend in data observations
 - Model fails to explain each data observation correctly (i.e. error)
- Let's recall our assumptions
 - There is a relationship between Olympics year and winning time
 - This relationship is linear
 - This relationship will hold in future
- Are these good assumptions?
- Still, ignoring the error is not right
 - Let's consider error as noise and model it

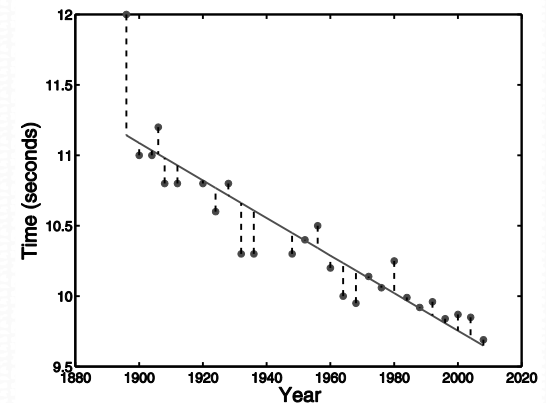
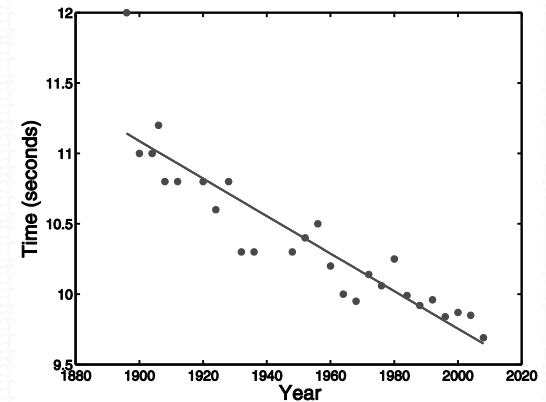


Linear modelling

- Approaches to linear modelling
 - Minimize the loss: minimize the squared error between model predictions and training observations
 - Maximize the likelihood: maximize the likelihood of match between model predictions and training observations
 - Build a model to generate data
 - Explicitly model the noise (i.e. the error between model and observations)

Thinking generatively

- Let's recall that our assumptions are weak
 - The process generating this data is very complex
- Still, can we try to build a model that can *generate* such data?
- *Generative modelling*: build a model to *generate* data that *looks like* data observations
 - $t_n = \mathbf{w}^T \mathbf{x}_n$
 - Error?
 - Modelling error



Refresher:

Probability

- Random variable (X, Y)
- Probability of tossing a coin (head=1, tail=0)
 - $P(X = 0) = 0.5$
 - $P(X = 1) = 0.5$

Refresher:

Conditional probability

- When the outcome of an event is affected (i.e. conditioned) by the outcome of another event
- For example: we toss a coin and then tell the result
 - $P(X = 1)$: probability of coin landing head
 - $P(X = 0)$: probability of coin landing tail
 - $P(Y = 1)$: probability of telling coin landed head
 - $P(Y = 0)$: probability of telling coin landed tail
- $P(Y = y|X = x)$: probability of telling coin landed y , given that coin has landed x

Refresher:

Conditional probability

- If we always tell the true outcome:
 - $P(Y = 1|X = 1) = ?$
 - $P(Y = 1|X = 1) = 1$
 - $P(Y = 0|X = 0) = ?$
 - $P(Y = 0|X = 0) = 1$
 - $P(Y = 0|X = 1) = ?$
 - $P(Y = 0|X = 1) = 0$
 - $P(Y = 1|X = 0) = ?$
 - $P(Y = 1|X = 0) = 0$

Refresher:

Conditional probability

- If we tell the true head outcome only 80% times:
 - $P(Y = 1|X = 1) = ?$
 - $P(Y = 1|X = 1) = 0.8$

 - $P(Y = 0|X = 0) = ?$
 - $P(Y = 0|X = 0) = 1$

 - $P(Y = 0|X = 1) = ?$
 - $P(Y = 0|X = 1) = 0.2$

 - $P(Y = 1|X = 0) = ?$
 - $P(Y = 1|X = 0) = 0$

Refresher:

Joint probability

- What is the probability that the coin lands heads and we say heads?
 - This is joint probability (i.e. probability of two or more variables)
 - $P(Y = y, X = x)$
- In case of no dependence between variables:
 - $P(Y = y, X = x) = P(Y = y)P(X = x)$
- In case of dependence between variables:
 - $P(Y = y, X = x) = P(Y = y|X = x)P(X = x)$, or
 - $P(Y = y, X = x) = P(X = x|Y = y)P(Y = y)$

Refresher:

Joint probability

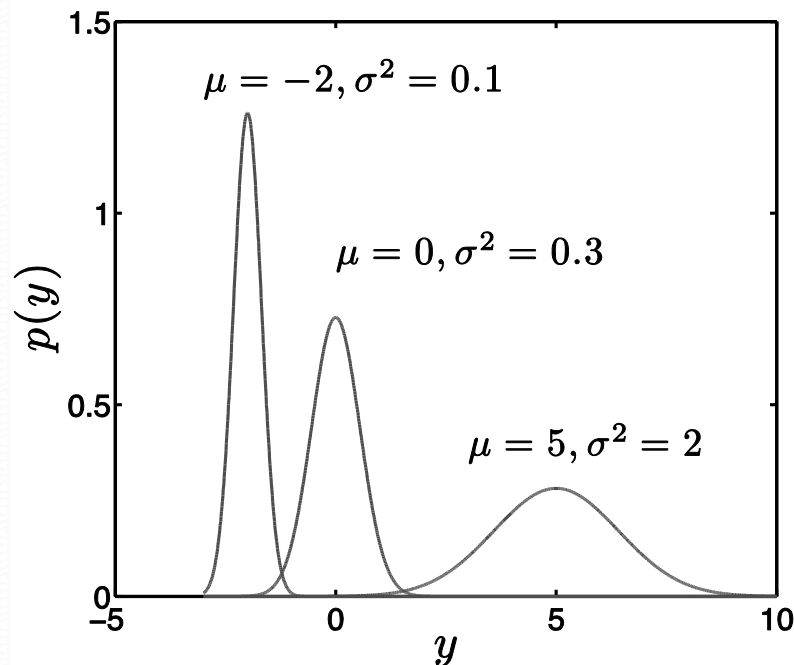
- What is the probability that the coin lands heads and we say heads?
 - $P(Y = 1, X = 1) = P(Y = 1|X = 1)P(X = 1) = ?$
 - $P(Y = 1, X = 1) = P(Y = 1|X = 1)P(X = 1) = 0.8 \times 0.5 = 0.4$
- Other joint probabilities
 - $P(Y = 0, X = 1) = P(Y = 0|X = 1)P(X = 1) = ?$
 - $P(Y = 0, X = 1) = P(Y = 0|X = 1)P(X = 1) = 0.2 \times 0.5 = 0.1$
 - $P(Y = 1, X = 0) = P(Y = 1|X = 0)P(X = 0) = ?$
 - $P(Y = 1, X = 0) = P(Y = 1|X = 0)P(X = 0) = 0 \times 0.5 = 0$
 - $P(Y = 0, X = 0) = P(Y = 0|X = 0)P(X = 0) = ?$
 - $P(Y = 0, X = 0) = P(Y = 0|X = 0)P(X = 0) = 1 \times 0.5 = 0.5$

Refresher:

Gaussian pdf

- Gaussian (or normal) probability density function

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



Thinking generatively

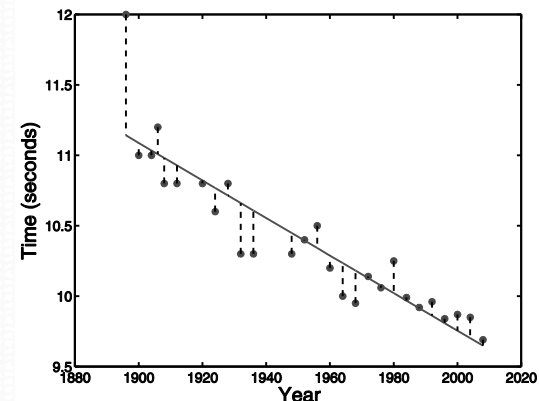
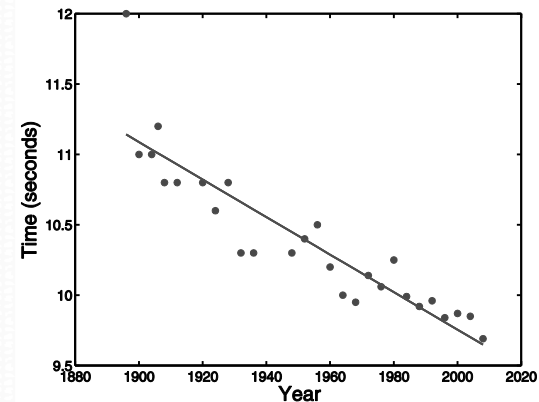
- *Generative modelling*: build a model to *generate* data that *looks like* data observations

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

- Generative model with noise considerations:

- $t_n = \mathbf{w}^T \mathbf{x}_n + \varepsilon_n$
- Additive noise?

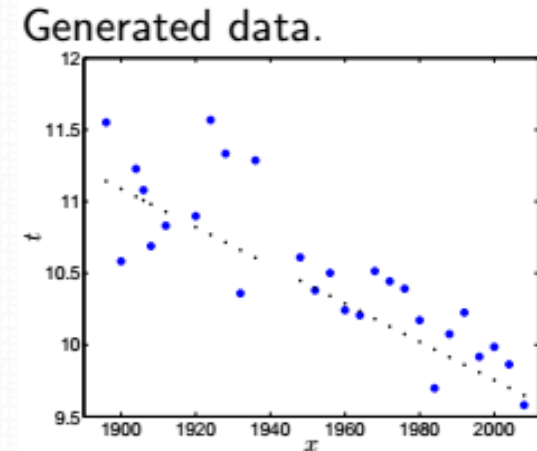
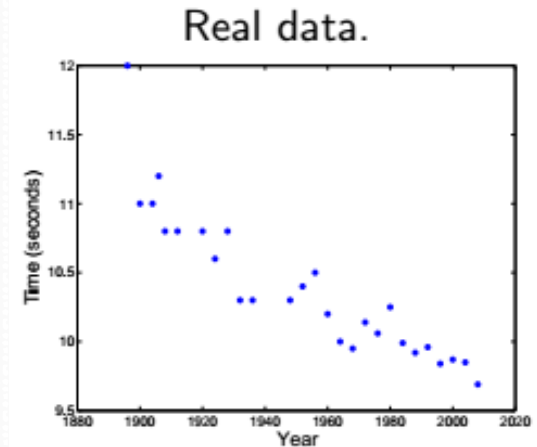
- Noise (ε_n)
 - Is it a random variable?
 - Is it discrete or continuous?
 - What pdf for ε_n ?



Thinking generatively

- Probability density of ε_n
$$p(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N) = \prod_{n=1}^N p(\varepsilon_n)$$
i.e. each ε_n is independent
- Let's use a Gaussian density for $p(\varepsilon_n) = \mathcal{N}(\mu, \sigma^2) = \mathcal{N}(0, 0.05)$

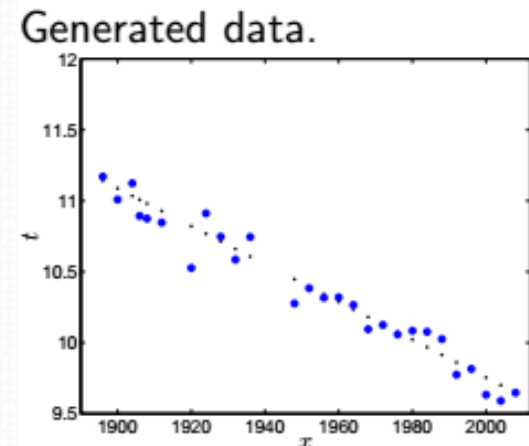
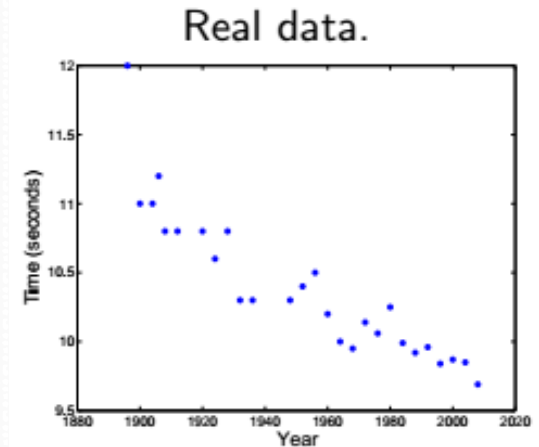
- Generative model with noise considerations:
 - $t_n = \mathbf{w}^T \mathbf{x}_n + \varepsilon_n$
- Deterministic component (i.e. trend)
- Random component (i.e. noise)



Thinking generatively

- Probability density of ε_n
$$p(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N) = \prod_{n=1}^N p(\varepsilon_n)$$
i.e. each ε_n is independent
- Let's use a Gaussian density for $p(\varepsilon_n) = \mathcal{N}(\mu, \sigma^2) = \mathcal{N}(0, 0.01)$

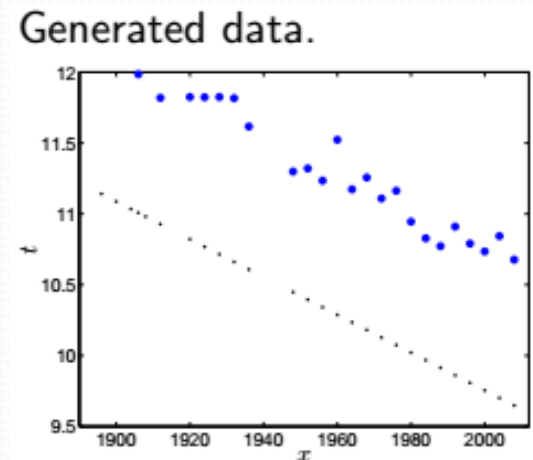
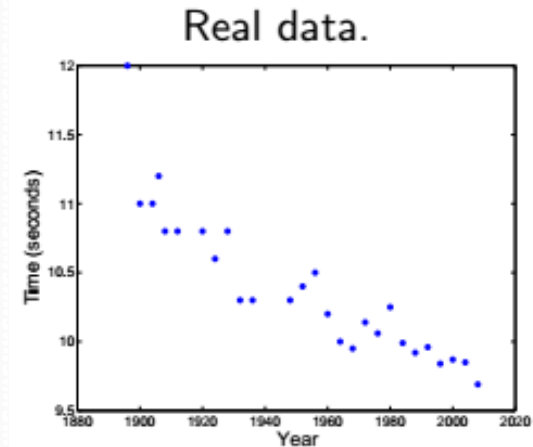
- Generative model with noise considerations:
 - $t_n = \mathbf{w}^T \mathbf{x}_n + \varepsilon_n$
- Deterministic component (i.e. trend)
- Random component (i.e. noise)



Thinking generatively

- Probability density of ε_n
$$p(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N) = \prod_{n=1}^N p(\varepsilon_n)$$
i.e. each ε_n is independent
- Let's use a Gaussian density for $p(\varepsilon_n) = \mathcal{N}(\mu, \sigma^2) = \mathcal{N}(1, 0.01)$

- Generative model with noise considerations:
 - $t_n = \mathbf{w}^T \mathbf{x}_n + \varepsilon_n$
- Deterministic component (i.e. trend)
- Random component (i.e. noise)



Thinking generatively

- Generative model with noise considerations:
 - $t_n = \mathbf{w}^T \mathbf{x}_n + \varepsilon_n$
 - $t_n = f(\mathbf{x}_n; \mathbf{w}) + \varepsilon_n$, where $p(\varepsilon_n) = \mathcal{N}(0, \sigma^2)$
 - Model is now determined not only by \mathbf{w} but also σ^2
- t_n can now be considered a random variable itself, due to addition of ε_n
 - i.e. for a given \mathbf{x}_n , t_n is not a single fixed value but rather is drawn out from a pdf
 - Thus finding \mathbf{w} and σ^2 by minimizing loss (with least squares approach) is not possible

Thinking generatively

- The probability density of t_n is:

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

- Let's recall $p(\varepsilon_n) = \mathcal{N}(0, \sigma^2)$

$$y = a + z$$

$$p(z) = \mathcal{N}(m, s)$$

$$p(y) = \mathcal{N}(m + a, s)$$

- Note that $\mathbf{w}^T \mathbf{x}_n$ determines the mean (i.e. trend) and σ^2 determines variance (i.e. noise)

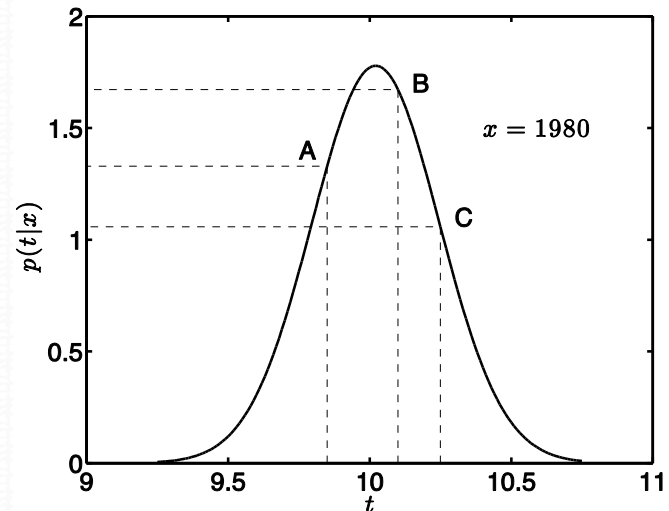
Thinking generatively

- Let's use this generative model to look at pdf for year 1980

$$p(t_n | \mathbf{x}_n = [1, 1980]^T, \mathbf{w} = [36.416, -0.0133]^T, \sigma^2 = 0.05)$$

This generates a Gaussian pdf $\mathcal{N}(\mathbf{w}^T \mathbf{x}_n = 10.02, \sigma^2 = 0.05)$

- Height of the curve corresponds to how likely it is to observe a particular t
- A, B, or C – more likely?
- Likelihood at $t_n = 10.25$?
- Can we determine \mathbf{w} and σ^2 that maximize likelihood of observed data t_n with generated data?



Likelihood

- For each input-response pair (\mathbf{x}_n, t_n) , we have a Gaussian likelihood: $p(t_n|\mathbf{x}_n, \mathbf{w}, \sigma^2)$
- Maximize the likelihood of matching all observed responses (t_1, t_2, \dots, t_N) conditioned on the observed data $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ and model parameters (\mathbf{w}, σ^2)

$$p(t_1, t_2, \dots, t_N | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N, \mathbf{w}, \sigma^2) = p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \sigma^2)$$

- Let's recall that $p(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N) = \prod_{n=1}^N p(\varepsilon_n)$

i.e. each ε_n is independent

- Thus, likelihood can be estimated as:

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

Maximum likelihood

- Likelihood estimate

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

- How to find \mathbf{w} and σ^2 ?
- Find “best” \mathbf{w} and σ^2 that maximize likelihood L

$$\underset{\mathbf{w}, \sigma^2}{\operatorname{argmax}} L$$

- It's mathematically convenient if we, instead, maximize the log of likelihood L

$$\underset{\mathbf{w}, \sigma^2}{\operatorname{argmax}} \log(L)$$

- Model parameters (\mathbf{w}, σ^2) that maximize log likelihood ($\log(L)$) also maximize likelihood (L)

Maximum likelihood

- Likelihood estimate

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

- Log likelihood estimate

$$\log(L) = \log \left[\prod_{n=1}^N \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2) \right] = \sum_{n=1}^N \log[\mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2)]$$

- Let's recall

$$\mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y - \mu)^2 \right\}$$

so

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

Maximum likelihood

$$\log(L) = \sum_{n=1}^N \log[\mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2)]$$

- Considering that

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

we get

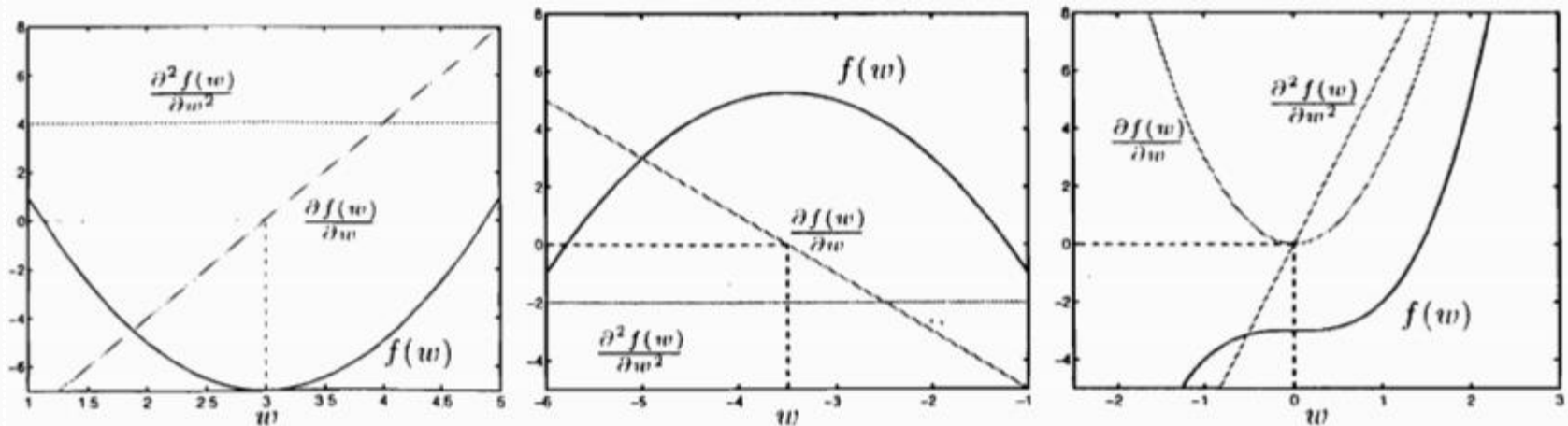
$$\log(L) = \sum_{n=1}^N \log\left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} (t_n - \mathbf{w}^T \mathbf{x}_n)^2\right\}\right]$$

which can be simplified to:

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

Finding function's maximum

- A function's minimum or maximum can be determined where the 1st derivative is zero
 - Local maxima
 - Local minima



Maximum likelihood

- Finding function's maximum

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

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

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

$$\frac{\partial \log(L)}{\partial \mathbf{w}} = \frac{1}{\sigma^2} \sum_{n=1}^N (\mathbf{x}_n t_n - \mathbf{x}_n \mathbf{x}_n^T \mathbf{w}) = \mathbf{0}$$

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

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

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

- It's exactly the same solution as from least squares
 - Minimizing squared loss = maximizing likelihood

Maximum likelihood

- Finding function's maximum

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

$$\frac{\partial \log(L)}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{n=1}^N (t_n - \mathbf{x}_n^T \hat{\mathbf{w}})^2 = 0$$

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

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

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

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix}$$

Model complexity

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

- Consider that:

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

so the log likelihood estimate at maximum is:

$$\log(L) = -\frac{N}{2}\log(2\pi) - N\log(\sigma) - \frac{1}{2\widehat{\sigma^2}} N\widehat{\sigma^2}$$

$$\log(L) = -\frac{N}{2}(1 + \log(2\pi)) - \frac{N}{2}\log(\widehat{\sigma^2})$$

Model complexity

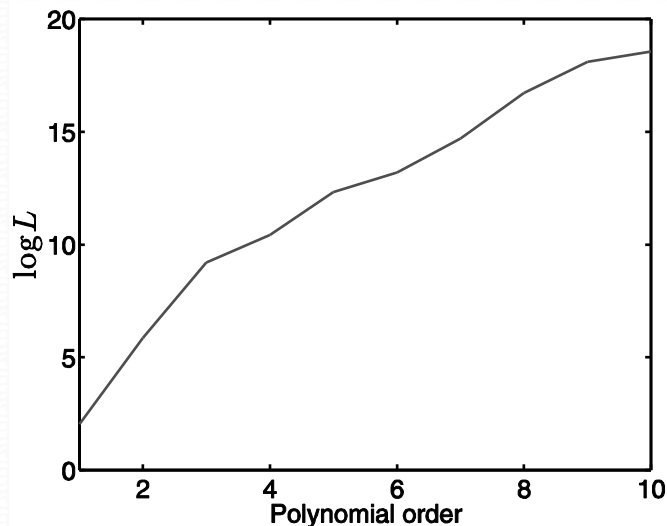
- Log likelihood estimate at maximum:

$$\log(L) = -\frac{N}{2}(1 + \log(2\pi)) - \frac{N}{2}\log(\widehat{\sigma^2})$$

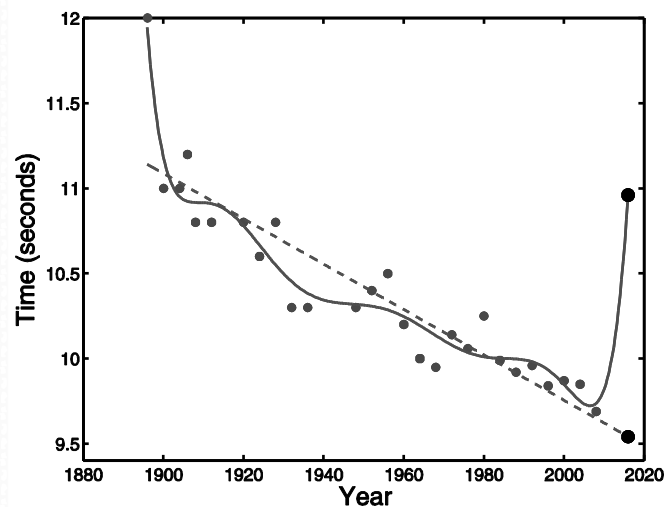
- Decrease in $\widehat{\sigma^2}$ will result in increase in log likelihood
- Note that $\mathbf{w}^T \mathbf{x}_n$ determines the mean (i.e. trend) and σ^2 determines variance (i.e. noise)
- Decrease in σ^2 will result in those values of model parameters \mathbf{w} that capture observations closely
 - i.e. over-fitting and poor generalization

Model complexity

- Modelling Olympics data again..
- Higher model complexity results in maximum likelihood
 - Generalization and over-fitting tradeoff



(a) Increase in log likelihood as the polynomial order increases



(b) 1st and 8th order polynomial functions fitted to the Olympics men's 100 m data. Large dark circles correspond to predictions for the 2016 Olympics

Bias-variance tradeoff

- Generalization and over-fitting tradeoff
- Error between predicted values and observed values can be considered to consist of two components – bias and model variance:

$$\bar{\mathcal{M}} = \mathcal{B}^2 + \mathcal{V}$$

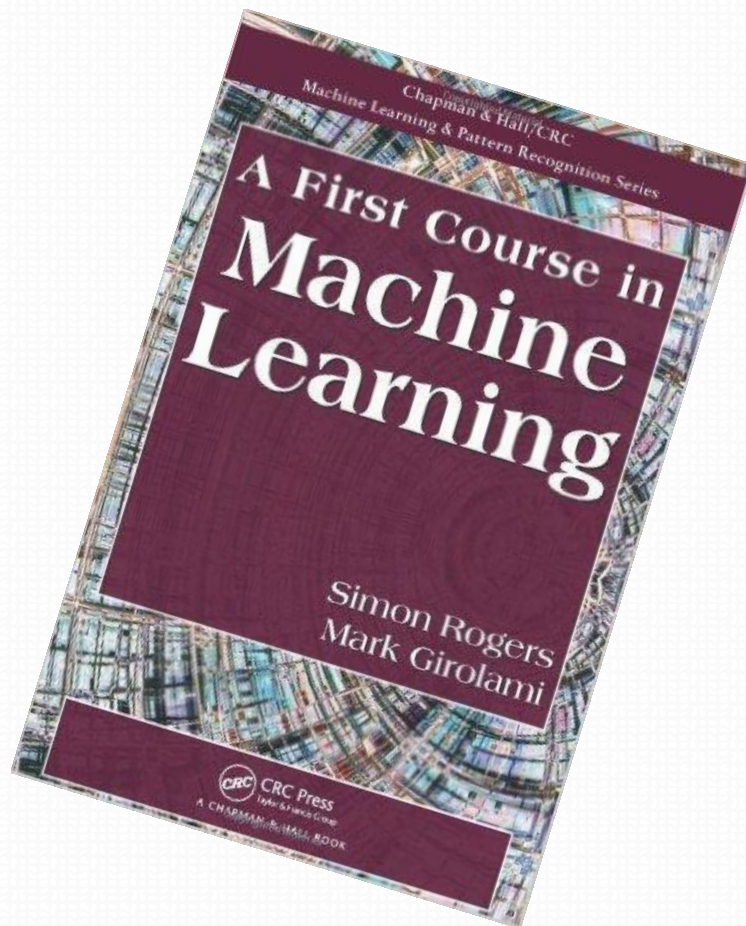
- \mathcal{B}^2 – bias: systematic mismatch between our model and the actual process that generated data
 - Decrease in bias \mathcal{B}^2 to control $\bar{\mathcal{M}}$
 - Complex models lead to decrease in bias
- \mathcal{V} – variance: more complex model has higher variance

Summary

- Explicit modelling of error as noise
- Noise modelling as Gaussian random variable
- Likelihood of model predictions and observed data
- Maximizing the likelihood
- Generalization and over-fitting tradeoff

Exercise (ungraded)

- Book (FCML) – exercise 2.1
- Book (FCML) – exercise 2.2 (refresher: probability)
- Book (FCML) – exercise 2.5 (refresher: probability)
- Book (FCML) – MATLAB code – `olymplike.m`
- Book (FCML) – MATLAB code – `genolymp.m`



Author's material
(Simon Rogers)



Thank You