## Statistical Machine Learning

Assignment Project Exam Help

Christian Walder + Lexing Xie

https://powcoder.com

Machine Learning Research Group, CSIRO Data61
ANU Computer Science

Canberra
Add WeChat powcoder
Semester One, 2021.

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")

Part I

*Linear Regression 2*
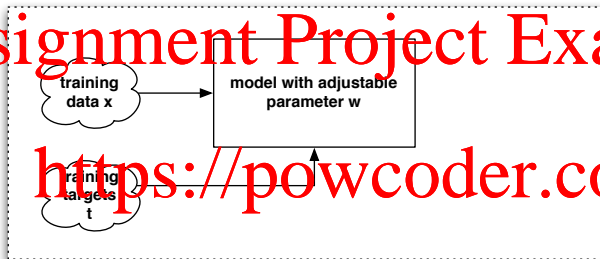
Assignment Project Exam Help

- Basis functions
- Maximum Likelihood with Gaussian Noise
- Regularisation
- Bias variance decomposition

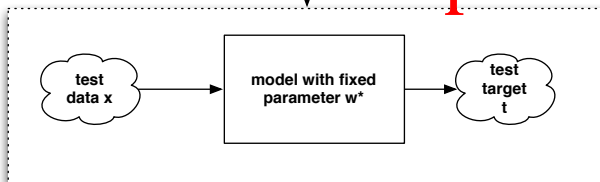https://powcoder.com

Add WeChat powcoder

# *Training and Testing:*
# *(Non-Bayesian) Point Estimate*

*Bayesian Regression*

Statistical Machine
Learning

©2021
Ong & Walder & Webers
& Xie
Data61 \ CSIRO
ANU Computer Science

- Bayes' theorem

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{normalisation}} \qquad p(\mathbf{w} \,|\, \mathbf{t}) = \frac{p(\mathbf{t} \,|\, \mathbf{w}) \, p(\mathbf{w})}{p(\mathbf{t})}$$

where we left out the conditioning on $\mathbf{x}$ (always assumed), and $\beta$, which is assumed to be constant.

- I.i.d. regression likelihood for additive Gaussian noise is

$$p(\mathbf{t} \,|\, \mathbf{w}) = \prod_{n=1}^{N} \mathcal{N}(t_n \,|\, y(\mathbf{x}_n, \mathbf{w}), \beta^{-1})$$

$$= \prod_{n=1}^{N} \mathcal{N}(t_n \,|\, \mathbf{w}^\top \phi(\mathbf{x}_n), \beta^{-1})$$

$$= \text{const} \times \exp\{-\beta \frac{1}{2} (\mathbf{t} - \mathbf{\Phi}\mathbf{w})^\top (\mathbf{t} - \mathbf{\Phi}\mathbf{w})\}$$

$$= \mathcal{N}(\mathbf{t} \,|\, \mathbf{\Phi}\mathbf{w}, \beta^{-1}\mathbf{I})$$

- The choice of prior affords an intuitive control over our inductive bias.
- All inference schemes have such biases, and often arise more opaquely than the prior in Bayes' rule.
- Can we find a prior for the given likelihood which
  - makes sense for the problem at hand
  - allows us to find a posterior in a 'nice' form

An answer to the second question:

**Definition (Conjugate Prior)**

A class of prior probability distributions $p(w)$ is conjugate to a class of likelihood functions $p(x \mid w)$ if the resulting posterior distributions $p(w \mid x)$ are in the same family as $p(w)$.

*Table:* Discrete likelihood distributions

| Likelihood | Conjugate Prior |
|------------|-----------------|
| Bernoulli | Beta |
| Binomial | Beta |
| Poisson | Gamma |
| Multinomial | Dirichlet |

*Table:* Continuous likelihood distributions

| Likelihood | Conjugate Prior |
|------------|-----------------|
| Uniform | Pareto |
| Exponential | Gamma |
| Normal | Normal (mean parameter) |
| Multivariate normal | Multivariate normal (mean parameter) |

- Example : If the likelihood function is Gaussian, choosing a Gaussian prior for the mean will ensure that the posterior distribution is also Gaussian.

- Given a marginal distribution for $\mathbf{x}$ and a conditional Gaussian distribution for $\mathbf{y}$ given $\mathbf{x}$ in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y} \mid \boldsymbol{A}\mathbf{x} + \mathbf{b}, \boldsymbol{L}^{-1})$$

- we get

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} \mid \boldsymbol{A}\boldsymbol{\mu} + \mathbf{b}, \boldsymbol{L}^{-1} + \boldsymbol{A}\boldsymbol{\Lambda}^{-1}\boldsymbol{A}^{\top})$$
$$p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\Sigma}\{\boldsymbol{A}^{\top}\boldsymbol{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$

where $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \boldsymbol{A}^{\top}\boldsymbol{L}\boldsymbol{A})^{-1}$.

Note that the covariance $\boldsymbol{\Sigma}$ does not involve $y$.

# *Conjugate Prior to a Gaussian Distribution (intuition)*

Given

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \,|\, \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\mathbf{y} \,|\, \mathbf{x}) = \mathcal{N}(\mathbf{y} \,|\, A\mathbf{x} + \mathbf{b}, L^{-1}) \Leftrightarrow \mathbf{y} = A\mathbf{x} + \mathbf{b} + \mathcal{N}(\mathbf{0}, L^{-1})$$

We have $\mathbb{E}[\mathbf{y}] = \mathbb{E}[A\mathbf{x} + \mathbf{b}] = A\boldsymbol{\mu} + \mathbf{b}$ and by the easily proven Bienaymé formula for the variance of the sum of uncorrelated variables,

$$\text{cov}[\mathbf{y}] = \underbrace{\text{cov}[A\mathbf{x} + \mathbf{b}]}_{\mathbb{E}[A\mathbf{x}(A\mathbf{x})^\top] = A\mathbb{E}[\mathbf{x}\mathbf{x}^\top]A^\top = A\boldsymbol{\Lambda}^{-1}A^\top} + \underbrace{\text{cov}[\mathcal{N}(\mathbf{0}, L^{-1})]}_{L^{-1}}.$$

So $\mathbf{y}$ is Gaussian with

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} \,|\, A\boldsymbol{\mu} + \mathbf{b}, L^{-1} + A\boldsymbol{\Lambda}^{-1}A^\top)$$

Then letting $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + A^\top L A)^{-1}$ and

$$p(\mathbf{x} \,|\, \mathbf{y}) = \mathcal{N}(\mathbf{x} \,|\, \boldsymbol{\Sigma}\{A^\top L(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$

$$\Leftrightarrow \mathbf{x} = \boldsymbol{\Sigma}\{A^\top L(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\} + \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$$

yields the correct moments for $\mathbf{x}$, since

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}[\boldsymbol{\Sigma}\{A^\top L(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}] = \boldsymbol{\Sigma}\{A^\top L(A\boldsymbol{\mu} + \mathbf{b} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}$$

$$= \boldsymbol{\Sigma}\{A^\top L A\boldsymbol{\mu} + \boldsymbol{\Lambda}\boldsymbol{\mu}\} = (\boldsymbol{\Lambda} + A^\top L A)^{-1}\{A^\top L A + \boldsymbol{\Lambda}\}\boldsymbol{\mu} = \boldsymbol{\mu},$$

and it is similar (but tedious ; don't do it) to recover $\text{cov}[\mathbf{x}] = \boldsymbol{\Lambda}$.

Statistical Machine
Learning
© 2021
Ong & Walder & Webers
& Xie
Data61 | CSIRO
ANU Computer Science

# *Bayesian Regression*

- Choose a Gaussian prior with mean $\mathbf{m}_0$ and covariance $\mathbf{S}_0$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_0, \mathbf{S}_0)$$

- Same likelihood as before (here written in vector form):

$$p(\mathbf{t} \mid \mathbf{w}, \beta) = \mathcal{N}(\mathbf{t} \mid \mathbf{\Phi}\mathbf{w}, \beta^{-1}I)$$

- Given $N$ data pairs $(\mathbf{x}_n, t_n)$, the posterior is

$$p(\mathbf{w} \mid \mathbf{t}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N)$$

where

$$\mathbf{m}_N = \mathbf{S}_N(\mathbf{S}_0^{-1}\mathbf{m}_0 + \beta\mathbf{\Phi}^\top\mathbf{t})$$
$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta\mathbf{\Phi}^\top\mathbf{\Phi}$$

(derive this with the identities on the previous slides)

# *Bayesian Regression: Zero Mean, Isotropic Prior*

- For simplicity we proceed with $\mathbf{m}_0 = 0$ and $\mathbf{S}_0 = \alpha^{-1}\mathbf{I}$, so

$$p(\mathbf{w} \mid \alpha) = \mathcal{N}(\mathbf{w} \mid 0, \alpha^{-1}\mathbf{I})$$

- The posterior becomes $p(\mathbf{w} \mid \mathbf{t}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N)$ with

$$\mathbf{m}_N = \beta \mathbf{S}_N \mathbf{\Phi}^\top \mathbf{t}$$
$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \mathbf{\Phi}^\top \mathbf{\Phi}$$

- For $\alpha \ll \beta$ we get

$$\mathbf{m}_N \to \mathbf{w}_{ML} = (\mathbf{\Phi}^\top \mathbf{\Phi})^{-1} \mathbf{\Phi}^\top \mathbf{t}$$

- Log of posterior is sum of log likelihood and log of prior

$$\ln p(\mathbf{w} \mid \mathbf{t}) = -\frac{\beta}{2}(\mathbf{t} - \mathbf{\Phi}\mathbf{w})^\top(\mathbf{t} - \mathbf{\Phi}\mathbf{w}) - \frac{\alpha}{2}\mathbf{w}^\top \mathbf{w} + \text{const}$$

# *Bayesian Regression*

- Log of posterior is sum of log likelihood and log of prior

$$\ln p(\mathbf{w}|\mathbf{t}) = -\beta \underbrace{\frac{1}{2}\|\mathbf{t} - \mathbf{\Phi}\mathbf{w}\|^2}_{\text{sum-of-squares-error}} - \underbrace{\frac{\alpha}{2}\|\mathbf{w}\|^2}_{\text{regulariser}} + \text{const.}$$

- The *maximum a posteriori* estimator

$$\mathbf{w}_{\text{m.a.p.}} = \arg\max_{\mathbf{w}} p(\mathbf{w}|\mathbf{t})$$

corresponds to minimising the sum-of-squares error function with quadratic regularisation coefficient $\lambda = \alpha/\beta$.

- The posterior is Gaussian so mode = mean: $\mathbf{w}_{\text{m.a.p.}} = \mathbf{m}_N$.

- For $\alpha \ll \beta$ the we recover unregularised least squares (equivalently m.a.p. approaches maximum likelihood), for example in case of
  - an infinitely broad prior with $\alpha \to 0$
  - an infinitely precise likelihood with $\beta \to \infty$

- If we have not yet seen any data point ($N = 0$), the posterior is equal to the prior.



- Sequential arrival of data points : the posterior given some observed data acts as the prior for the future data.
- Nicely fits a sequential learning framework.

- Example of a linear basis function model
- Single input $x$, single output $t$
- Linear model $y(x, \mathbf{w}) = w_0 + w_1 x$.
- True data distribution / sampling procedure:
  1. Choose an $x_n$ from the uniform distribution $\mathcal{U}(x \mid -1, +1)$.
  2. Calculate $f(x_n, \mathbf{a}) = a_0 + a_1 x_n$, where $a_0 = -0.3$, $a_1 = 0.5$.
  3. Add Gaussian noise with standard deviation $\sigma = 0.2$,
  $$t_n \sim \mathcal{N}(x_n, f(x_n, \mathbf{a}), 0.04)$$
- Set the precision of the uniform prior to $\alpha = 2.0$.

- In the training phase, data **x** and targets **t** are provided
- In the test phase, a new data value $x$ is given and the corresponding target value $t$ is asked for
- Bayesian approach: Find the probability of the test target $t$ given the test data $x$, the training data **x** and the training targets **t**

$$p(t \,|\, x, \mathbf{x}, \mathbf{t})$$

- This is the Predictive Distribution (c.f. the posterior distribution, which is over the parameters).

# *How to calculate the Predictive Distribution?*

- Introduce the model parameter $\mathbf{w}$ via the sum rule

$$p(t \,|\, x, \mathbf{x}, \mathbf{t}) = \int p(t, \mathbf{w} \,|\, x, \mathbf{x}, \mathbf{t})d\mathbf{w}$$
$$= \int p(t \,|\, \mathbf{w}, x, \mathbf{x}, \mathbf{t})p(\mathbf{w} \,|\, x, \mathbf{x}, \mathbf{t})d\mathbf{w}$$

- The test target $t$ depends only on the test data $x$ and the model parameter $\mathbf{w}$, but not on the training data and the training targets

$$p(t \,|\, \mathbf{w}, x, \mathbf{x}, \mathbf{t}) = p(t \,|\, \mathbf{w}, x)$$

- The model parameter $\mathbf{w}$ were learned with the training data $\mathbf{x}$ and the training targets $\mathbf{t}$ only

$$p(\mathbf{w} \,|\, x, \mathbf{x}, \mathbf{t}) = p(\mathbf{w} \,|\, \mathbf{x}, \mathbf{t})$$

- Predictive Distribution

$$p(t \,|\, x, \mathbf{x}, \mathbf{t}) = \int p(t \,|\, \mathbf{w}, x)p(\mathbf{w} \,|\, \mathbf{x}, \mathbf{t})d\mathbf{w}$$

The predictive distribution is

$$p(t \,|\, x, \mathbf{x}, \mathbf{t}) = \int p(t \,|\, \mathbf{w}, x, \mathbf{x}, \mathbf{t}) p(\mathbf{w} \,|\, x, \mathbf{x}, \mathbf{t}) d\mathbf{w}$$

because

$$\int p(t \,|\, \mathbf{w}, x, \mathbf{x}, \mathbf{t}) p(\mathbf{w} \,|\, x, \mathbf{x}, \mathbf{t}) d\mathbf{w} = \int \frac{p(t, \mathbf{w}, x, \mathbf{x}, \mathbf{t})}{p(\mathbf{w}, x, \mathbf{x}, \mathbf{t})} \frac{p(\mathbf{w}, x, \mathbf{x}, \mathbf{t})}{p(x, \mathbf{x}, \mathbf{t})} d\mathbf{w}$$

$$= \int \frac{p(t, \mathbf{w}, x, \mathbf{x}, \mathbf{t})}{p(x, \mathbf{x}, \mathbf{t})} d\mathbf{w}$$

$$= \frac{p(t, x, \mathbf{x}, \mathbf{t})}{p(x, \mathbf{x}, \mathbf{t})}$$

$$= p(t \,|\, x, \mathbf{x}, \mathbf{t}),$$

or simply

$$\int p(t \,|\, \mathbf{w}, x, \mathbf{x}, \mathbf{t}) p(\mathbf{w} \,|\, x, \mathbf{x}, \mathbf{t}) d\mathbf{w} = \int p(t, \mathbf{w} \,|\, x, \mathbf{x}, \mathbf{t}) d\mathbf{w}$$

$$= p(t \,|\, x, \mathbf{x}, \mathbf{t}).$$

*Statistical Machine Learning*

© 2021
Ong & Walder & Webers
& Xie
*Data61 | CSIRO*
*ANU Computer Science*

# *Predictive Distribution with Simplified Prior*

- Find the predictive distribution

$$p(t \mid \mathbf{t}, \alpha, \beta) = \int p(t \mid \mathbf{w}, \beta) p(\mathbf{w} \mid \mathbf{t}, \alpha, \beta) d\mathbf{w}$$

(remember : conditioning on $\mathbf{x}$ is often suppressed to simplify the notation.)

- Now we know (neglecting as usual to notate conditioning on $\mathbf{x}$)

$$p(t \mid \mathbf{w}, \beta) = \mathcal{N}(t \mid \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}), \beta^{-1})$$

- and the posterior was

$$p(\mathbf{w} \mid \mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N)$$

where

$$\mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^\top \mathbf{t}$$
$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^\top \boldsymbol{\Phi}$$

# *Predictive Distribution with Simplified Prior*

- If we do the integral (it turns out to be the convolution of the two Gaussians), we get for the predictive distribution

$$p(t \mid \mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t \mid \mathbf{m}_N^\top \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

  where the variance $\sigma_N^2(\mathbf{x})$ is given by

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^\top \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}).$$

- This is more easily shown using a similar approach to the earlier "intuition" slide and again with the Bienaymé formula, now using

$$t = \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x}) + \mathcal{N}(0, \beta^{-1}).$$

  However this is a linear-Gaussian specific trick and in general we need to integrate out the parameters.

Example with artificial sinusoidal data from $\sin(2\pi x)$ (green) and added noise. Number of data points $N = 1$.



Mean of the predictive distribution (red) and regions of one standard deviation from mean (red shaded).

Example with artificial sinusoidal data from $\sin(2\pi x)$ (green) and added noise. Number of data points $N = 2$.



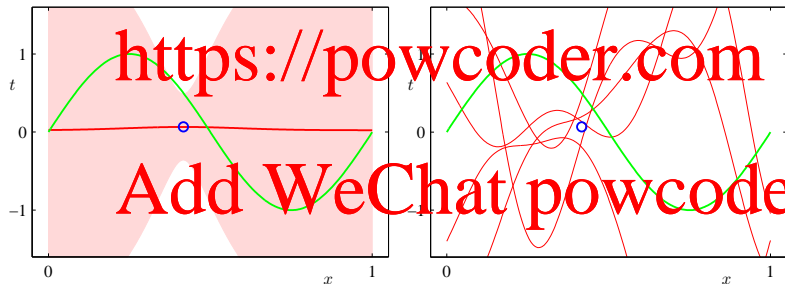Mean of the predictive distribution (red) and regions of one standard deviation from mean (red shaded).

Example with artificial sinusoidal data from $\sin(2\pi x)$ (green) and added noise. Number of data points $N = 4$.



Mean of the predictive distribution (red) and regions of one standard deviation from mean (red shaded).

Example with artificial sinusoidal data from $\sin(2\pi x)$ (green)
and added noise. Number of data points $N = 25$.



Mean of the predictive distribution (red) and regions of one
standard deviation from mean (red shaded).

Plots of the function $y(x, \mathbf{w})$ using samples from the posterior distribution over $\mathbf{w}$. Number of data points $N = 1$.
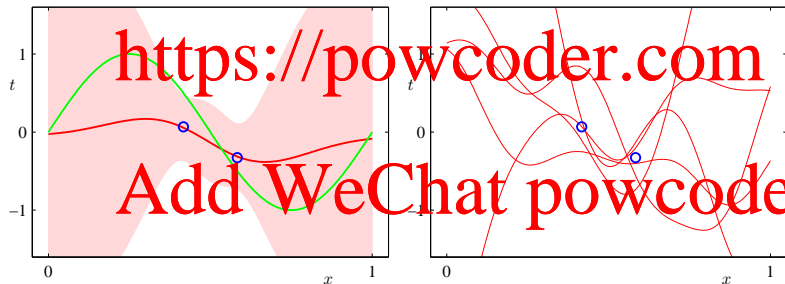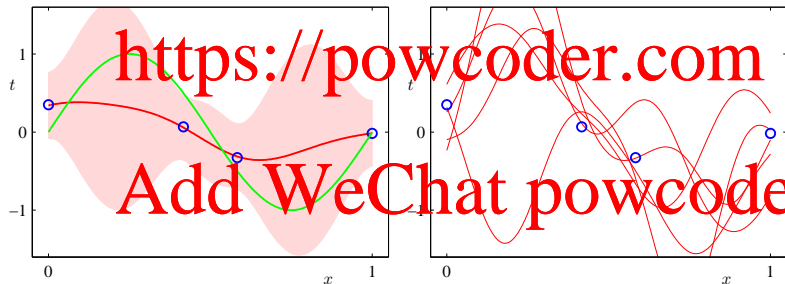
Plots of the function $y(x, \mathbf{w})$ using samples from the posterior distribution over $\mathbf{w}$. Number of data points $N = 2$.

Plots of the function $y(x, \mathbf{w})$ using samples from the posterior distribution over $\mathbf{w}$. Number of data points $N = 4$.

Plots of the function $y(x, \mathbf{w})$ using samples from the posterior distribution over $\mathbf{w}$. Number of data points $N = 25$.

- Basis function $\phi_i(\mathbf{x})$ are fixed before the training data set is observed.
- Curse of dimensionality : Number of basis function grows rapidly, often exponentially, with the dimensionality $D$.
- But typical data sets have two nice properties which can be exploited if the basis functions are not fixed :
  - Data lie close to a nonlinear manifold with intrinsic dimension much smaller than $D$. Need algorithms which place basis functions only where data are (*e.g.* kernel methods / Gaussian processes).
  - Target variables may only depend on a few significant directions within the data manifold. Need algorithms which can exploit this property (*e.g.* linear methods or shallow neural networks).

- Linear Algebra allows us to operate in $n$-dimensional vector spaces using the intuition from our 3-dimensional world as a vector space. No surprises as long as $n$ is finite.

- If we add more structure to a vector space (e.g. inner product, metric), our intuition gained from the 3-dimensional world around us may be wrong.

- Example: Sphere of radius $r = 1$. What is the fraction of the volume of the sphere in a $D$-dimensional space which lies between radius $r = 1$ and $r = 1 - \epsilon$ ?

- Volume scales like $r^D$, therefore the formula for the volume of a sphere is $V_D(r) = K_D r^D$.

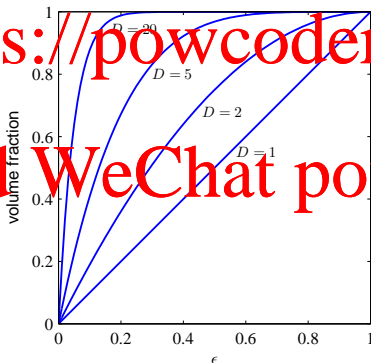$$\frac{V_D(1) - V_D(1 - \epsilon)}{V_D(1)} = 1 - (1 - \epsilon)^D$$

- Fraction of the volume of the sphere in a $D$-dimensional space which lies between radius $r = 1$ and $r = 1 - \epsilon$

$$\frac{V_D(1) - V_D(1-\epsilon)}{V_D(1)} = 1 - (1-\epsilon)^D$$
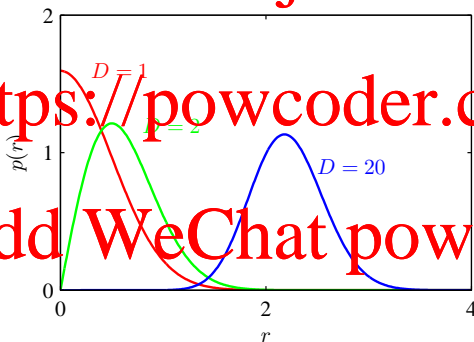
- Probability density with respect to radius $r$ of a Gaussian distribution for various values of the dimensionality $D$.

# *Curse of Dimensionality*

- Probability density with respect to radius $r$ of a Gaussian distribution for various values of the dimensionality $D$.

- Example: $D = 2$; assume $\mu = 0, \Sigma = I$

$$\mathcal{N}(x \mid 0, I) = \frac{1}{2\pi} \exp\left\{-\frac{1}{2}x^\top x\right\} = \frac{1}{2\pi} \exp\left\{-\frac{1}{2}(x_1^2 + x_2^2)\right\}$$

- Coordinate transformation

$$x_1 = r\cos(\phi) \qquad x_2 = r\sin(\phi)$$

- Probability in the new coordinates

$$p(r, \phi \mid 0, I) = \mathcal{N}(r(x), \phi(x) \mid 0, I) \mid J \mid$$

where $|J| = r$ is the determinant of the Jacobian for the given coordinate transformation.

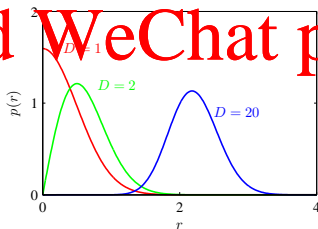$$p(r, \phi \mid 0, I) = \frac{1}{2\pi} r \exp\left\{-\frac{1}{2}r^2\right\}$$

*Statistical Machine Learning*

© 2021
Ong & Walder & Webers
& Xie
Data61 \ CSIRO
ANU Computer Science

# Curse of Dimensionality

- Probability density with respect to radius $r$ of a Gaussian distribution for $D = 2$ (and $\mu = 0, \Sigma = I$)

$$p(r \mid 0, I) = \frac{1}{2\pi} r \exp\left\{-\frac{1}{2} r^2\right\}$$

- Integrate over all angles $\phi$

$$p(r \mid 0, I) = \int_0^{2\pi} \frac{1}{2\pi} r \exp\left\{-\frac{1}{2} r^2\right\} d\phi = r \exp\left\{-\frac{1}{2} r^2\right\}$$

- Basis functions
- Maximum likelihood with Gaussian noise
- Regularisation
- Bias variance decomposition
- Conjugate prior
- Bayesian linear regression
- Sequential update of the posterior
- Predictive distribution
- Curse of dimensionality