# Linear Models for Regression
## Greg Mori - CMPT 419/726

Bishop PRML Ch. 3

# Outline

Regression

Linear Basis Function Models

Loss Functions for Regression

Finding Optimal Weights

Regularization

Bayesian Linear Regression

# Outline
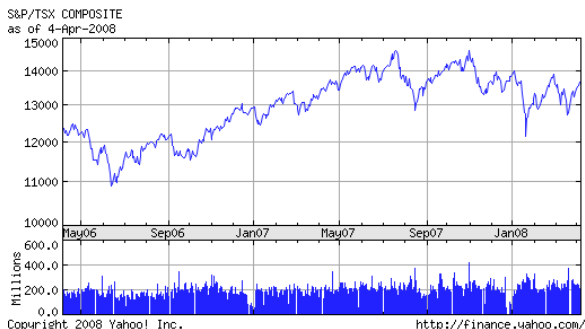
## Regression

Linear Basis Function Models

Loss Functions for Regression

Finding Optimal Weights

Regularization

Bayesian Linear Regression

# Regression



S&P/TSX COMPOSITE
as of 4-Apr-2008

Copyright 2008 Yahoo! Inc.                    http://finance.yahoo.com/

- Given training set $\{(\boldsymbol{x}_1, \boldsymbol{t}_1), \ldots, (\boldsymbol{x}_N, \boldsymbol{t}_N)\}$
- $\boldsymbol{t}_i$ is continuous: regression
- For now, assume $t_i \in \mathbb{R}$, $\boldsymbol{x}_i \in \mathbb{R}^D$
- E.g. $\boldsymbol{t}_i$ is stock price, $\boldsymbol{x}_i$ contains company profit, debt, cash flow, gross sales, number of spam emails sent, . . .

# Outline

# Linear Functions

- A function $f(\cdot)$ is linear if

$$f(\alpha u + \beta v) = \alpha f(u) + \beta f(v)$$

- Linear functions will lead to simple algorithms, so let's see what we can do with them

# Linear Regression

- Simplest linear model for regression

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 x_1 + w_2 x_2 + \ldots + w_D x_D$$

  - Remember, we're learning $\boldsymbol{w}$
  - Set $\boldsymbol{w}$ so that $y(\boldsymbol{x}, \boldsymbol{w})$ aligns with target value in training data

# Linear Regression

- Simplest linear model for regression

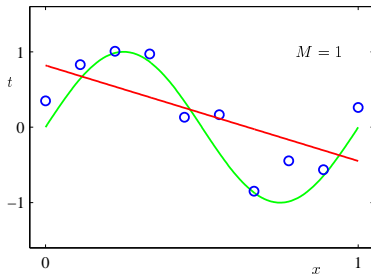$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 x_1 + w_2 x_2 + \ldots + w_D x_D$$

  - Remember, we're learning $\boldsymbol{w}$
  - Set $\boldsymbol{w}$ so that $y(\boldsymbol{x}, \boldsymbol{w})$ aligns with target value in training data

# Linear Regression

- Simplest linear model for regression

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 x_1 + w_2 x_2 + \ldots + w_D x_D$$

  - Remember, we're learning $\boldsymbol{w}$
  - Set $\boldsymbol{w}$ so that $y(\boldsymbol{x}, \boldsymbol{w})$ aligns with target value in training data
- This is a very simple model, limited in what it can do

# Linear Basis Function Models

- Simplest linear model

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 x_1 + w_2 x_2 + \ldots + w_D x_D$$

  was linear in $\boldsymbol{x}$ ($^*$) and $\boldsymbol{w}$

- Linear in $\boldsymbol{w}$ is what will be important for simple algorithms

- Extend to include fixed non-linear functions of data

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 \phi_1(\boldsymbol{x}) + w_2 \phi_2(\boldsymbol{x}) + \ldots + w_{M-1} \phi_{M-1}(\boldsymbol{x})$$

- Linear combinations of these basis functions also linear in parameters

# Linear Basis Function Models

- Simplest linear model

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 x_1 + w_2 x_2 + \ldots + w_D x_D$$

  was linear in $\boldsymbol{x}$ (*) and $\boldsymbol{w}$

- Linear in $\boldsymbol{w}$ is what will be important for simple algorithms

- Extend to include fixed non-linear functions of data

  $$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 \phi_1(\boldsymbol{x}) + w_2 \phi_2(\boldsymbol{x}) + \ldots + w_{M-1} \phi_{M-1}(\boldsymbol{x})$$

- Linear combinations of these basis functions also linear in parameters

# Linear Basis Function Models

- Simplest linear model

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 x_1 + w_2 x_2 + \ldots + w_D x_D$$

  was linear in $\boldsymbol{x}$ (*) and $\boldsymbol{w}$

- Linear in $\boldsymbol{w}$ is what will be important for simple algorithms
- Extend to include fixed non-linear functions of data

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1 \phi_1(\boldsymbol{x}) + w_2 \phi_2(\boldsymbol{x}) + \ldots + w_{M-1} \phi_{M-1}(\boldsymbol{x})$$

- Linear combinations of these basis functions also linear in parameters

# Linear Basis Function Models

- Bias parameter allows fixed offset in data

$$y(\boldsymbol{x}, \boldsymbol{w}) = \underbrace{w_0}_{bias} + w_1\phi_1(\boldsymbol{x}) + w_2\phi_2(\boldsymbol{x}) + \ldots + w_{M-1}\phi_{M-1}(\boldsymbol{x})$$

- Think of simple 1-D $x$:

$$y(x, w) = \underbrace{w_0}_{intercept} + \underbrace{w_1}_{slope} x$$

- For notational convenience, define $\phi_0(x) = 1$:

$$y(\boldsymbol{x}, \boldsymbol{w}) = \sum_{j=0}^{M-1} w_j\phi_j(\boldsymbol{x}) = \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x})$$

# Linear Basis Function Models

- Bias parameter allows fixed offset in data

$$y(\boldsymbol{x}, \boldsymbol{w}) = \underbrace{w_0}_{bias} + w_1\phi_1(\boldsymbol{x}) + w_2\phi_2(\boldsymbol{x}) + \ldots + w_{M-1}\phi_{M-1}(\boldsymbol{x})$$

- Think of simple 1-D $\boldsymbol{x}$:

$$y(x, w) = \underbrace{w_0}_{intercept} + \underbrace{w_1}_{slope} x$$

- For notational convenience, define $\phi_0(x) = 1$:

$$y(\boldsymbol{x}, \boldsymbol{w}) = \sum_{j=0}^{M-1} w_j\phi_j(\boldsymbol{x}) = \boldsymbol{w}^T\phi(\boldsymbol{x})$$

# Linear Basis Function Models

- Bias parameter allows fixed offset in data

$$y(\boldsymbol{x}, \boldsymbol{w}) = \underbrace{w_0}_{bias} + w_1\phi_1(\boldsymbol{x}) + w_2\phi_2(\boldsymbol{x}) + \ldots + w_{M-1}\phi_{M-1}(\boldsymbol{x})$$

  - Think of simple 1-D $x$:

$$y(x, w) = \underbrace{w_0}_{intercept} + \underbrace{w_1}_{slope} x$$

- For notational convenience, define $\phi_0(x) = 1$:

$$y(\boldsymbol{x}, \boldsymbol{w}) = \sum_{j=0}^{M-1} w_j\phi_j(\boldsymbol{x}) = \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x})$$

# Linear Basis Function Models

- Function for regression $y(\boldsymbol{x}, \boldsymbol{w})$ is non-linear function of $\boldsymbol{x}$, but linear in $\boldsymbol{w}$:

$$y(\boldsymbol{x}, \boldsymbol{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x})$$

- Polynomial regression is an example of this
- Order $M$ polynomial regression, $\phi_j(x) = ?$

# Linear Basis Function Models

- Function for regression $y(\boldsymbol{x}, \boldsymbol{w})$ is non-linear function of $\boldsymbol{x}$, but linear in $\boldsymbol{w}$:

$$y(\boldsymbol{x}, \boldsymbol{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x})$$

- Polynomial regression is an example of this
- Order $M$ polynomial regression, $\phi_j(x) = ?$

# Linear Basis Function Models

- Function for regression $y(\boldsymbol{x}, \boldsymbol{w})$ is non-linear function of $\boldsymbol{x}$, but linear in $\boldsymbol{w}$:

$$y(\boldsymbol{x}, \boldsymbol{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x})$$

- Polynomial regression is an example of this
- Order $M$ polynomial regression, $\phi_j(x) = ?$
- $\phi_j(x) = x^j$:

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 x^0 + w_1 x^1 + \ldots + w_M x^M$$

## Basis Functions: Feature Functions

- Often we extract features from $x$
    - An intuitve way to think of $\phi_j(x)$ is as feature functions
- E.g. Automatic CMPT726 project report grading system
    - $x$ is text of report:  In this project we apply the
      algorithm of Mori [2] to recognizing blue
      objects.  We test this algorithm on
      pictures of you and I from my holiday photo
      collection.  ...
- $\phi_1(x)$ is count of occurrences of Mori [
- $\phi_2(x)$ is count of occurrences of of you and I
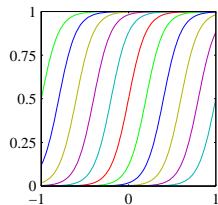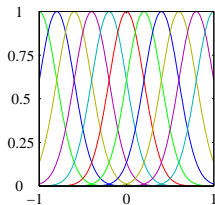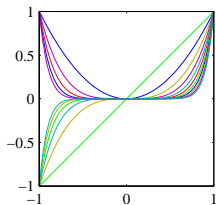- Regression grade $y(x, w) = 20\phi_1(x) - 10\phi_2(x)$

## Basis Functions: Feature Functions

- Often we extract features from $x$
    - An intuitve way to think of $\phi_j(x)$ is as feature functions
- E.g. Automatic CMPT726 project report grading system
    - $x$ is text of report:      In this project we apply the algorithm of Mori [2] to recognizing blue objects.  We test this algorithm on pictures of you and I from my holiday photo collection.  ...
- $\phi_1(x)$ is count of occurrences of Mori [
- $\phi_2(x)$ is count of occurrences of of you and I
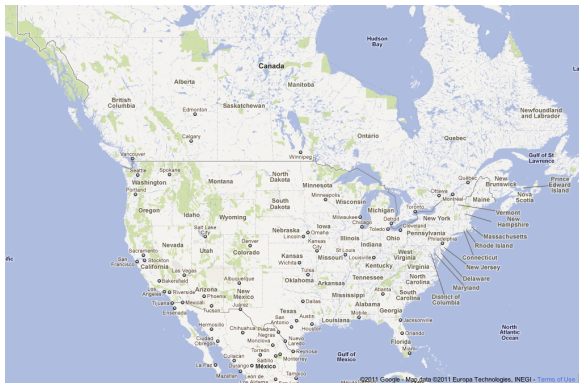- Regression grade $y(x, w) = 20\phi_1(x) - 10\phi_2(x)$

## Basis Functions: Feature Functions

- Often we extract features from $x$
  - An intuitve way to think of $\phi_j(x)$ is as feature functions
- E.g. Automatic CMPT726 project report grading system
  - $x$ is text of report:  In this project we apply the algorithm of Mori [2] to recognizing blue objects. We test this algorithm on pictures of you and I from my holiday photo collection. ...
- $\phi_1(x)$ is count of occurrences of Mori [
- $\phi_2(x)$ is count of occurrences of of you and I
- Regression grade $y(x, w) = 20\phi_1(x) - 10\phi_2(x)$

## Basis Functions: Feature Functions

- Often we extract features from $x$
  - An intuitve way to think of $\phi_j(x)$ is as feature functions
- E.g. Automatic CMPT726 project report grading system
  - $x$ is text of report:      In this project we apply the
    algorithm of Mori [2] to recognizing blue
    objects.  We test this algorithm on
    pictures of you and I from my holiday photo
    collection.   ...
- $\phi_1(x)$ is count of occurrences of Mori [
- $\phi_2(x)$ is count of occurrences of of you and I
- Regression grade $y(x, w) = 20\phi_1(x) - 10\phi_2(x)$

# Other Non-linear Basis Functions
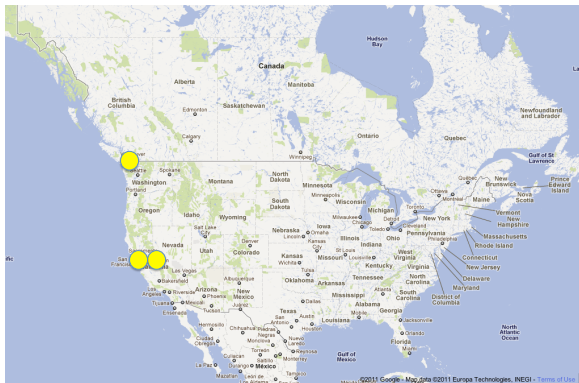


- Polynomial $\phi_j(x) = x^j$
- Gaussians $\phi_j(x) = \exp\{-\frac{(x-\mu_j)^2}{2s^2}\}$
- Sigmoidal $\phi_j(x) = \frac{1}{1+\exp((\mu_j-x)/s)}$
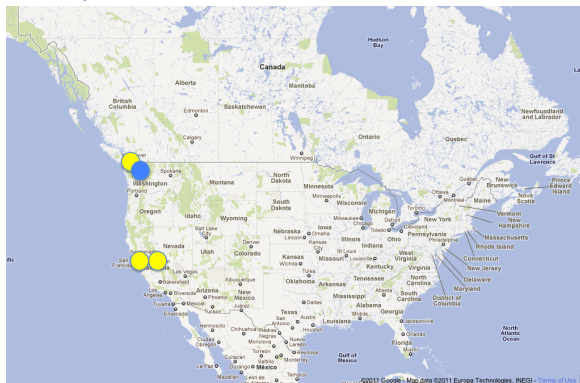
## Example - Gaussian Basis Functions: Temperature



- Use Gaussian basis functions, regression on temperature
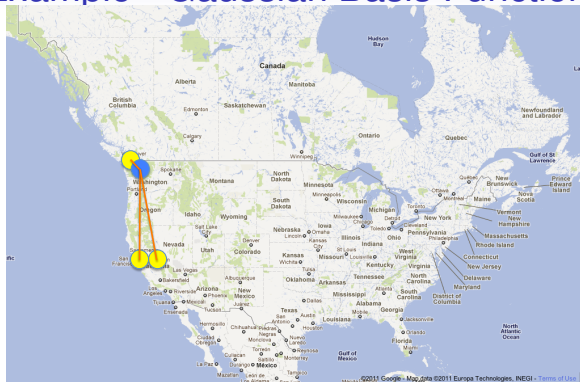
## Example - Gaussian Basis Functions: Temperature



- $\mu_1 = $ Vancouver, $\mu_2 = $ San Francisco, $\mu_3 = $ Oakland

# Example - Gaussian Basis Functions: Temperature



- $\mu_1 =$ Vancouver, $\mu_2 =$ San Francisco, $\mu_3 =$ Oakland
- Temperature in $x =$ Seattle? $y(\boldsymbol{x}, \boldsymbol{w}) =$
  $w_1 \exp\{-\frac{(x-\mu_1)^2}{2s^2}\} + w_2 \exp\{-\frac{(x-\mu_2)^2}{2s^2}\} + w_3 \exp\{-\frac{(x-\mu_3)^2}{2s^2}\}$

# Example - Gaussian Basis Functions: Temperature



- $\mu_1 = $ Vancouver, $\mu_2 = $ San Francisco, $\mu_3 = $ Oakland
- Temperature in $x = $ Seattle? $y(\boldsymbol{x}, \boldsymbol{w}) = $
  $w_1 \exp\{-\frac{(x-\mu_1)^2}{2s^2}\} + w_2 \exp\{-\frac{(x-\mu_2)^2}{2s^2}\} + w_3 \exp\{-\frac{(x-\mu_3)^2}{2s^2}\}$
- Compute distances to all $\mu$, $y(\boldsymbol{x}, \boldsymbol{w}) \approx w_1$

# Example - Gaussian Basis Functions: 726 Report Grading

- Define:
  - $\mu_1 =$ Crime and Punishment
  - $\mu_2 =$ Animal Farm
  - $\mu_3 =$ Some paper by Mori

- Learn weights:
  - $w_1 = ?$
  - $w_2 = ?$
  - $w_3 = ?$

- Grade a project report $x$:
  - Measure similarity of $x$ to each $\mu$, Gaussian, with weights:
    $y(x, w) =$
    $w_1 \exp\{-\frac{(x-\mu_1)^2}{2s^2}\} + w_2 \exp\{-\frac{(x-\mu_2)^2}{2s^2}\} + w_3 \exp\{-\frac{(x-\mu_3)^2}{2s^2}\}$

- The Gaussian basis function models end up similar to template matching

## Example - Gaussian Basis Functions: 726 Report Grading

- Define:
    - $\mu_1 =$ Crime and Punishment
    - $\mu_2 =$ Animal Farm
    - $\mu_3 =$ Some paper by Mori
- Learn weights:
    - $w_1 = ?$
    - $w_2 = ?$
    - $w_3 = ?$
- Grade a project report $x$:
    - Measure similarity of $x$ to each $\mu$, Gaussian, with weights:
      $y(x, w) =$
      $w_1 \exp\{-\frac{(x-\mu_1)^2}{2s^2}\} + w_2 \exp\{-\frac{(x-\mu_2)^2}{2s^2}\} + w_3 \exp\{-\frac{(x-\mu_3)^2}{2s^2}\}$
- The Gaussian basis function models end up similar to template matching

## Example - Gaussian Basis Functions: 726 Report Grading

- Define:
  - $\mu_1 =$ Crime and Punishment
  - $\mu_2 =$ Animal Farm
  - $\mu_3 =$ Some paper by Mori
- Learn weights:
  - $w_1 = ?$
  - $w_2 = ?$
  - $w_3 = ?$
- Grade a project report $x$:
  - Measure similarity of $x$ to each $\mu$, Gaussian, with weights:
    $y(x, w) =$
    $w_1 \exp\{-\frac{(x-\mu_1)^2}{2s^2}\} + w_2 \exp\{-\frac{(x-\mu_2)^2}{2s^2}\} + w_3 \exp\{-\frac{(x-\mu_3)^2}{2s^2}\}$
- The Gaussian basis function models end up similar to template matching

# Outline

## Loss Functions for Regression

- We want to find the "best" set of coefficients $w$
- Recall, one way to define "best" was minimizing squared error:

$$E(w) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, w) - t_n\}^2$$

- We will now look at another way, based on maximum likelihood

## Gaussian Noise Model for Regression

- We are provided with a training set $\{(\boldsymbol{x}_i, t_i)\}$
- Let's assume $t$ arises from a deterministic function plus Gassian distributed (with precision $\beta$) noise:

$$t = y(\boldsymbol{x}, \boldsymbol{w}) + \epsilon$$

- The probability of observing a target value $t$ is then:

$$p(t|\boldsymbol{x}, \boldsymbol{w}, \beta) = \mathcal{N}(t|y(\boldsymbol{x}, \boldsymbol{w}), \beta^{-1})$$

- Notation: $\mathcal{N}(x|\mu, \sigma^2)$; $x$ drawn from Gaussian with mean $\mu$, variance $\sigma^2$

## Gaussian Noise Model for Regression

- We are provided with a training set $\{(\boldsymbol{x}_i, t_i)\}$
- Let's assume $t$ arises from a deterministic function plus Gassian distributed (with precision $\beta$) noise:

$$t = y(\boldsymbol{x}, \boldsymbol{w}) + \epsilon$$

- The probability of observing a target value $t$ is then:

$$p(t|\boldsymbol{x}, \boldsymbol{w}, \beta) = \mathcal{N}(t|y(\boldsymbol{x}, \boldsymbol{w}), \beta^{-1})$$

- Notation: $\mathcal{N}(x|\mu, \sigma^2)$; $x$ drawn from Gaussian with mean $\mu$, variance $\sigma^2$

## Maximum Likelihood for Regression

- The likelihood of data $\boldsymbol{t} = \{t_i\}$ using this Gaussian noise model is:

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

- The log-likelihood is:

$$\ln p(\boldsymbol{t}|\boldsymbol{w}, \beta) = \ln \prod_{n=1}^{N} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp(-\frac{\beta}{2}(t_n - \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x}_n))^2)$$

$$= \underbrace{\frac{N}{2}\ln\beta - \frac{N}{2}\ln(2\pi)}_{const.\ wrt\ w} - \beta \underbrace{\frac{1}{2}\sum_{n=1}^{N}(t_n - \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x}_n))^2}_{squared\ error}$$

- Sum of squared errors is maximum likelihood under a Gaussian noise model

## Maximum Likelihood for Regression

- The likelihood of data $t = \{t_i\}$ using this Gaussian noise model is:

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

- The log-likelihood is:

$$\ln p(\boldsymbol{t}|\boldsymbol{w}, \beta) = \ln \prod_{n=1}^{N} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp(-\frac{\beta}{2}(t_n - \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x}_n))^2)$$

$$= \underbrace{\frac{N}{2}\ln\beta - \frac{N}{2}\ln(2\pi)}_{const.\ wrt\ w} - \beta \underbrace{\frac{1}{2}\sum_{n=1}^{N}(t_n - \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x}_n))^2}_{squared\ error}$$

- Sum of squared errors is maximum likelihood under a Gaussian noise model

## Maximum Likelihood for Regression

- The likelihood of data $t = \{t_i\}$ using this Gaussian noise model is:

$$p(t|w, \beta) = \prod_{n=1}^{N} \mathcal{N}(t_n|w^T\phi(x_n), \beta^{-1})$$

- The log-likelihood is:

$$\ln p(t|w, \beta) = \ln \prod_{n=1}^{N} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp(-\frac{\beta}{2}(t_n - w^T\phi(x_n))^2)$$

$$= \underbrace{\frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)}_{const.\ wrt\ w} - \beta \underbrace{\frac{1}{2} \sum_{n=1}^{N} (t_n - w^T\phi(x_n))^2}_{squared\ error}$$

- Sum of squared errors is maximum likelihood under a Gaussian noise model

## Maximum Likelihood for Regression

- The likelihood of data $\boldsymbol{t} = \{t_i\}$ using this Gaussian noise model is:

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

- The log-likelihood is:

$$\ln p(\boldsymbol{t}|\boldsymbol{w}, \beta) = \ln \prod_{n=1}^{N} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp(-\frac{\beta}{2}(t_n - \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x}_n))^2)$$

$$= \underbrace{\frac{N}{2}\ln\beta - \frac{N}{2}\ln(2\pi)}_{const.\ wrt\ \boldsymbol{w}} - \beta \underbrace{\frac{1}{2}\sum_{n=1}^{N}(t_n - \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x}_n))^2}_{squared\ error}$$

- Sum of squared errors is maximum likelihood under a Gaussian noise model

# Outline

# Finding Optimal Weights

- How do we maximize likelihood wrt $w$ (or minimize squared error)?

- Take gradient of log-likelihood wrt $w$:

$$\frac{\partial}{\partial w_i} \ln p(t|w, \beta) = \beta \sum_{n=1}^{N} (t_n - w^T \phi(x_n))\phi_i(x_n)$$

- In vector form:

$$\nabla \ln p(t|w, \beta) = \beta \sum_{n=1}^{N} (t_n - w^T \phi(x_n))\phi(x_n)^T$$

# Finding Optimal Weights

- How do we maximize likelihood wrt $w$ (or minimize squared error)?
- Take gradient of log-likelihood wrt $w$:

$$\frac{\partial}{\partial w_i} \ln p(t|w, \beta) = \beta \sum_{n=1}^{N} (t_n - w^T \phi(x_n))\phi_i(x_n)$$

- In vector form:

$$\nabla \ln p(t|w, \beta) = \beta \sum_{n=1}^{N} (t_n - w^T \phi(x_n))\phi(x_n)^T$$

# Finding Optimal Weights

- How do we maximize likelihood wrt $w$ (or minimize squared error)?

- Take gradient of log-likelihood wrt $w$:

$$\frac{\partial}{\partial w_i} \ln p(t|w, \beta) = \beta \sum_{n=1}^{N} (t_n - w^T \phi(x_n)) \phi_i(x_n)$$

- In vector form:

$$\nabla \ln p(t|w, \beta) = \beta \sum_{n=1}^{N} (t_n - w^T \phi(x_n)) \phi(x_n)^T$$

# Finding Optimal Weights

- Set gradient to 0:

$$
\mathbf{0}^T = \sum_{n=1}^{N} t_n \phi(\boldsymbol{x}_n)^T - \boldsymbol{w}^T \left( \sum_{n=1}^{N} \phi(\boldsymbol{x}_n)\phi(\boldsymbol{x}_n)^T \right)
$$

- Maximum likelihood estimate for *w*:

$$
w_{ML} = \left( \Phi^T \Phi \right)^{-1} \Phi^T t
$$

$$
\Phi = \begin{pmatrix}
\phi_0(\boldsymbol{x}_1) & \phi_1(\boldsymbol{x}_1) & \dots & \phi_{M-1}(\boldsymbol{x}_1) \\
\phi_0(\boldsymbol{x}_2) & \phi_1(\boldsymbol{x}_2) & \dots & \phi_{M-1}(\boldsymbol{x}_2) \\
\vdots & \vdots & \ddots & \vdots \\
\phi_0(\boldsymbol{x}_N) & \phi_1(\boldsymbol{x}_N) & \dots & \phi_{M-1}(\boldsymbol{x}_N)
\end{pmatrix}
$$

- $\Phi^{\dagger} = \left( \Phi^T \Phi \right)^{-1} \Phi^T$ known as the pseudo-inverse (numpy.linalg.pinv in python)

# Finding Optimal Weights

- Set gradient to 0:

$$\mathbf{0}^T = \sum_{n=1}^{N} t_n \phi(\boldsymbol{x}_n)^T - \boldsymbol{w}^T \left( \sum_{n=1}^{N} \phi(\boldsymbol{x}_n) \phi(\boldsymbol{x}_n)^T \right)$$

- Maximum likelihood estimate for $\boldsymbol{w}$:

$$\boldsymbol{w}_{ML} = \left( \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\Phi}^T \boldsymbol{t}$$

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_0(\boldsymbol{x}_1) & \phi_1(\boldsymbol{x}_1) & \dots & \phi_{M-1}(\boldsymbol{x}_1) \\ \phi_0(\boldsymbol{x}_2) & \phi_1(\boldsymbol{x}_2) & \dots & \phi_{M-1}(\boldsymbol{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\boldsymbol{x}_N) & \phi_1(\boldsymbol{x}_N) & \dots & \phi_{M-1}(\boldsymbol{x}_N) \end{pmatrix}$$
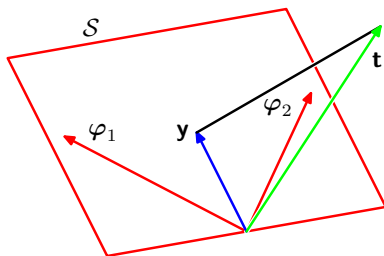
- $\boldsymbol{\Phi}^\dagger = \left( \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\Phi}^T$ known as the pseudo-inverse
  (`numpy.linalg.pinv` in python)

# Geometry of Least Squares



- $\boldsymbol{t} = (t_1, \ldots, t_N)$ is the target value vector
- $\mathcal{S}$ is space spanned by $\boldsymbol{\varphi}_j = (\phi_j(\boldsymbol{x}_1), \ldots, \phi_j(\boldsymbol{x}_N))$
- Solution $\boldsymbol{y}$ lies in $\mathcal{S}$
- Least squares solution is orthogonal projection of $\boldsymbol{t}$ onto $\mathcal{S}$
- Can verify this by looking at $\boldsymbol{y} = \Phi \boldsymbol{w}_{ML} = \Phi \Phi^\dagger \boldsymbol{t} = P \boldsymbol{t}$
  - $P^2 = P, P = P^T$

# Geometry of Least Squares



- $\boldsymbol{t} = (t_1, \ldots, t_N)$ is the target value vector
- $\mathcal{S}$ is space spanned by $\boldsymbol{\varphi}_j = (\phi_j(\boldsymbol{x}_1), \ldots, \phi_j(\boldsymbol{x}_N))$
- Solution $\boldsymbol{y}$ lies in $\mathcal{S}$
- Least squares solution is orthogonal projection of $\boldsymbol{t}$ onto $\mathcal{S}$
- Can verify this by looking at $\boldsymbol{y} = \Phi \boldsymbol{w}_{ML} = \Phi\Phi^\dagger t = Pt$
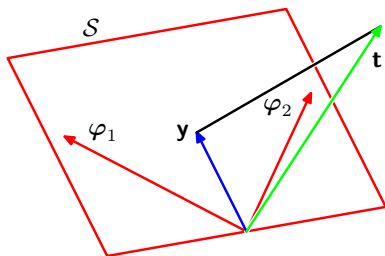  - $P^2 = P, P = P^T$

# Geometry of Least Squares



- $t = (t_1, \ldots, t_N)$ is the target value vector
- $\mathcal{S}$ is space spanned by $\varphi_j = (\phi_j(x_1), \ldots, \phi_j(x_N))$
- Solution $y$ lies in $\mathcal{S}$
- Least squares solution is orthogonal projection of $t$ onto $\mathcal{S}$
- Can verify this by looking at $y = \Phi w_{ML} = \Phi \Phi^\dagger t = Pt$
  - $P^2 = P, P = P^T$

# Geometry of Least Squares



- $t = (t_1, \ldots, t_N)$ is the target value vector
- $\mathcal{S}$ is space spanned by $\varphi_j = (\phi_j(x_1), \ldots, \phi_j(x_N))$
- Solution $y$ lies in $\mathcal{S}$
- Least squares solution is orthogonal projection of $t$ onto $\mathcal{S}$
- Can verify this by looking at $y = \Phi w_{ML} = \Phi\Phi^\dagger t = Pt$
  - $P^2 = P$, $P = P^T$

# Geometry of Least Squares



- $t = (t_1, \ldots, t_N)$ is the target value vector
- $\mathcal{S}$ is space spanned by $\varphi_j = (\phi_j(x_1), \ldots, \phi_j(x_N))$
- Solution $y$ lies in $\mathcal{S}$
- Least squares solution is orthogonal projection of $t$ onto $\mathcal{S}$
- Can verify this by looking at $y = \Phi w_{ML} = \Phi \Phi^{\dagger} t = Pt$
  - $P^2 = P, P = P^T$

# Sequential Learning

- In practice $N$ might be huge, or data might arrive online
- Can use a gradient descent method:
    - Start with initial guess for $w$
    - Update by taking a step in gradient direction $\nabla E$ of error function
- Modify to use stochastic / sequential gradient descent:
    - If error function $E = \sum_n E_n$ (e.g. least squares)
    - Update by taking a step in gradient direction $\nabla E_n$ for one example
    - Details about step size are important – decrease step size at the end

# Sequential Learning

- In practice $N$ might be huge, or data might arrive online
- Can use a gradient descent method:
    - Start with initial guess for $w$
    - Update by taking a step in gradient direction $\nabla E$ of error function
- Modify to use stochastic / sequential gradient descent:
    - If error function $E = \sum_n E_n$ (e.g. least squares)
    - Update by taking a step in gradient direction $\nabla E_n$ for one example
    - Details about step size are important – decrease step size at the end

# Sequential Learning

- In practice $N$ might be huge, or data might arrive online
- Can use a gradient descent method:
    - Start with initial guess for $w$
    - Update by taking a step in gradient direction $\nabla E$ of error function
- Modify to use stochastic / sequential gradient descent:
    - If error function $E = \sum_n E_n$ (e.g. least squares)
    - Update by taking a step in gradient direction $\nabla E_n$ for one example
    - Details about step size are important – decrease step size at the end

# Sequential Learning

- In practice $N$ might be huge, or data might arrive online
- Can use a gradient descent method:
    - Start with initial guess for $w$
    - Update by taking a step in gradient direction $\nabla E$ of error function
- Modify to use stochastic / sequential gradient descent:
    - If error function $E = \sum_n E_n$ (e.g. least squares)
    - Update by taking a step in gradient direction $\nabla E_n$ for one example
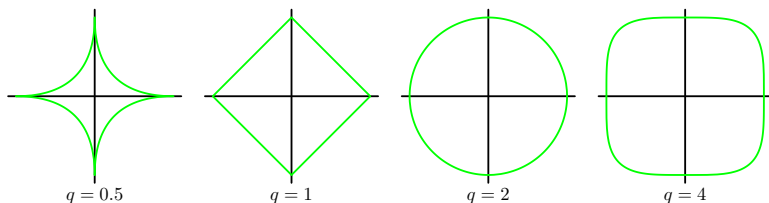    - Details about step size are important – decrease step size at the end

# Outline

# Regularization

- Last week we discussed regularization as a technique to avoid over-fitting:

$$\tilde{E}(\boldsymbol{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \boldsymbol{w}) - t_n\}^2 + \underbrace{\frac{\lambda}{2} ||\boldsymbol{w}||^2}_{regularizer}$$

- Next on the menu:
  - Other regularlizers
  - Bayesian learning and quadratic regularizer

# Other Regularizers



$q = 0.5$                $q = 1$                $q = 2$                $q = 4$

- Can use different norms for regularizer:

$$\tilde{E}(\boldsymbol{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \boldsymbol{w}) - t_n\}^2 + \frac{\lambda}{2} \sum_{j=1}^{M} |w_j|^q$$

- e.g. $q = 2$ – ridge regression
- e.g. $q = 1$ – lasso
- math is easiest with ridge regression

## Optimization with a Quadratic Regularizer

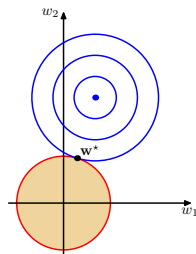- With $q = 2$, total error still a nice quadratic:

$$\tilde{E}(\boldsymbol{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \boldsymbol{w}) - t_n\}^2 + \frac{\lambda}{2} \boldsymbol{w}^T \boldsymbol{w}$$

- Calculus ...

$$\boldsymbol{w} = (\underbrace{\lambda \boldsymbol{I} + \boldsymbol{\Phi}^T \boldsymbol{\Phi}}_{regularized})^{-1} \boldsymbol{\Phi}^T \boldsymbol{t}$$

- Similar to unregularlized least squares
- Advantage $(\lambda \boldsymbol{I} + \boldsymbol{\Phi}^T \boldsymbol{\Phi})$ is well conditioned so inversion is stable

# Ridge Regression vs. Lasso



- Ridge regression aka parameter shrinkage
  - Weights $w$ shrink back towards origin

# Ridge Regression vs. Lasso



- Ridge regression aka parameter shrinkage
  - Weights $w$ shrink back towards origin
- Lasso leads to sparse models
  - Components of $w$ tend to 0 with large $\lambda$ (strong regularization)
  - Intuitively, once minimum achieved at large radius, minimum is on $w_1 = 0$

# Outline

# Bayesian Linear Regression

- Last week we saw an example of a Bayesian approach
  - Coin tossing - prior on parameters
- We will now do the same for linear regression
  - Prior on parameter $w$
- There will turn out to be a connection to regularlization

# Bayesian Linear Regression

- Start with a prior over parameters $w$
    - Conjugate prior is a Gaussian:

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

    - This simple form will make math easier; can be done for arbitrary Gaussian too

- Data likelihood, Gaussian model as before:

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

# Bayesian Linear Regression

- Start with a prior over parameters $w$
  - Conjugate prior is a Gaussian:

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

  - This simple form will make math easier; can be done for arbitrary Gaussian too

- Data likelihood, Gaussian model as before:

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

# Bayesian Linear Regression

- Posterior distribution on $w$:

$$p(w|t) \propto \left( \prod_{n=1}^{N} p(t_n|x_n, w, \beta) \right) p(w)$$

$$= \left[ \prod_{n=1}^{N} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left( -\frac{\beta}{2}(t_n - w^T \phi(x_n))^2 \right) \right] \left( \frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \exp(-\frac{\alpha}{2}w^T w)$$

- Take the log:

$$-\ln p(w|t) = \frac{\beta}{2} \sum_{n=1}^{N} (t_n - w^T \phi(x_n))^2 + \frac{\alpha}{2}w^T w + const$$

- $L_2$ regularization is maximum a posteriori (MAP) with a Gaussian prior.
  - $\lambda = \alpha/\beta$

# Bayesian Linear Regression

- Posterior distribution on $w$:

$$p(\boldsymbol{w}|\boldsymbol{t}) \propto \left( \prod_{n=1}^{N} p(t_n|\boldsymbol{x}_n, \boldsymbol{w}, \beta) \right) p(\boldsymbol{w})$$

$$= \left[ \prod_{n=1}^{N} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left( -\frac{\beta}{2}(t_n - \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x}_n))^2 \right) \right] \left( \frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \exp(-\frac{\alpha}{2}\boldsymbol{w}^T\boldsymbol{w})$$

- Take the log:

$$-\ln p(\boldsymbol{w}|\boldsymbol{t}) = \frac{\beta}{2} \sum_{n=1}^{N} (t_n - \boldsymbol{w}^T\boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\alpha}{2}\boldsymbol{w}^T\boldsymbol{w} + const$$

- $L_2$ regularization is maximum a posteriori (MAP) with a Gaussian prior.
  - $\lambda = \alpha/\beta$

# Bayesian Linear Regression

- Posterior distribution on $w$:

$$p(\boldsymbol{w}|\boldsymbol{t}) \propto \left( \prod_{n=1}^{N} p(t_n|\boldsymbol{x}_n, \boldsymbol{w}, \beta) \right) p(\boldsymbol{w})$$

$$= \left[ \prod_{n=1}^{N} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left( -\frac{\beta}{2}(t_n - \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_n))^2 \right) \right] \left( \frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \exp(-\frac{\alpha}{2}\boldsymbol{w}^T\boldsymbol{w})$$

- Take the log:

$$-\ln p(\boldsymbol{w}|\boldsymbol{t}) = \frac{\beta}{2} \sum_{n=1}^{N} (t_n - \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\alpha}{2} \boldsymbol{w}^T \boldsymbol{w} + const$$

- $L_2$ regularization is maximum a posteriori (MAP) with a Gaussian prior.
  - $\lambda = \alpha/\beta$

# Bayesian Linear Regression

- Posterior distribution on $w$:

$$p(\boldsymbol{w}|\boldsymbol{t}) \propto \left( \prod_{n=1}^{N} p(t_n|\boldsymbol{x}_n, \boldsymbol{w}, \beta) \right) p(\boldsymbol{w})$$

$$= \left[ \prod_{n=1}^{N} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left( -\frac{\beta}{2}(t_n - \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_n))^2 \right) \right] \left( \frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \exp(-\frac{\alpha}{2} \boldsymbol{w}^T \boldsymbol{w})$$

- Take the log:

$$-\ln p(\boldsymbol{w}|\boldsymbol{t}) = \frac{\beta}{2} \sum_{n=1}^{N} (t_n - \boldsymbol{w}^T \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\alpha}{2} \boldsymbol{w}^T \boldsymbol{w} + const$$
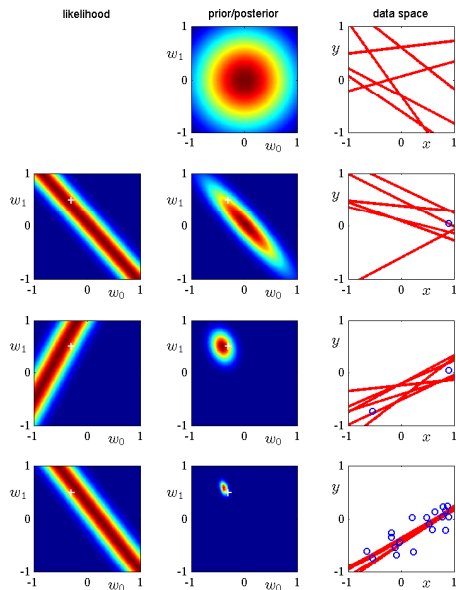
- $L_2$ regularization is maximum a posteriori (MAP) with a Gaussian prior.
  - $\lambda = \alpha/\beta$

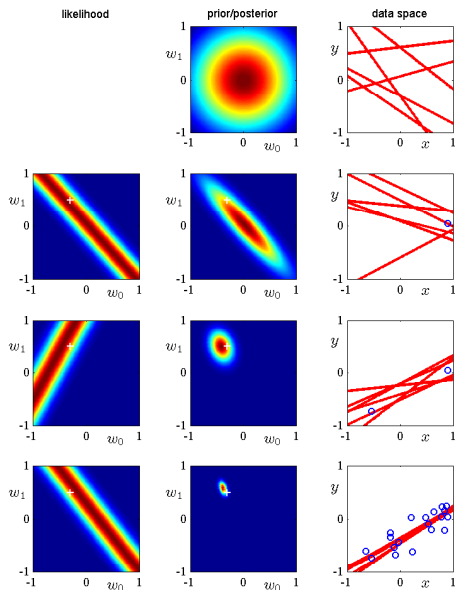# Bayesian Linear Regression - Intuition



- Simple example $x, t \in \mathbb{R}$, $y(x, \boldsymbol{w}) = w_0 + w_1 x$

- Start with Gaussian prior in parameter space

- Samples shown in data space

- Receive data points (blue circles in data space)

- Compute likelihood

- Posterior is prior (or prev. posterior) times likelihood

# Bayesian Linear Regression - Intuition



- Simple example $x, t \in \mathbb{R}$, $y(x, \boldsymbol{w}) = w_0 + w_1 x$
- Start with Gaussian prior in parameter space
- Samples shown in data space
- Receive data points (blue circles in data space)
- Compute likelihood
- Posterior is prior (or prev. posterior) times likelihood
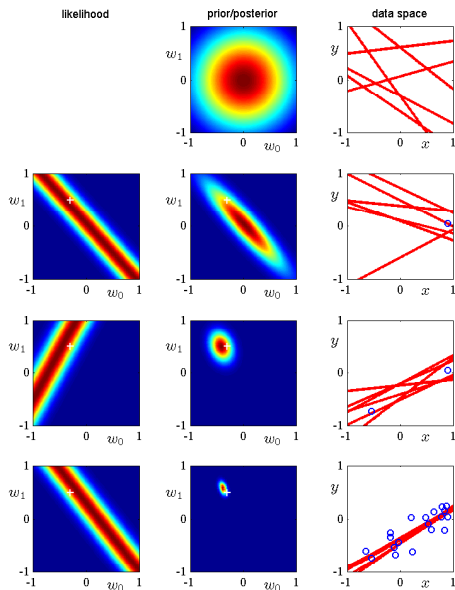
# Bayesian Linear Regression - Intuition



- Simple example $x, t \in \mathbb{R}$, $y(x, \boldsymbol{w}) = w_0 + w_1 x$
- Start with Gaussian prior in parameter space
- Samples shown in data space
- Receive data points (blue circles in data space)
- Compute likelihood
- Posterior is prior (or prev. posterior) times likelihood
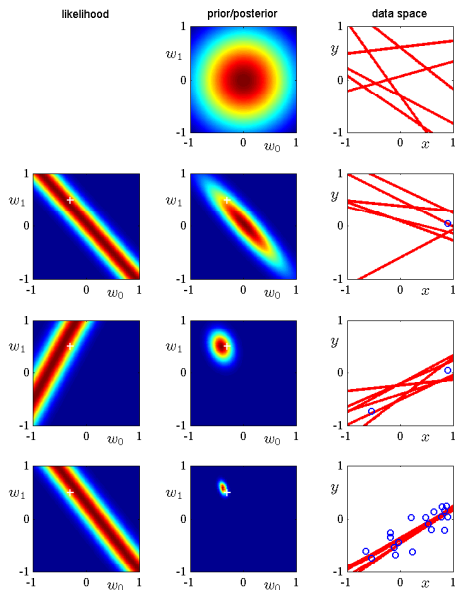
# Bayesian Linear Regression - Intuition



- Simple example $x, t \in \mathbb{R}$, $y(x, \boldsymbol{w}) = w_0 + w_1 x$
- Start with Gaussian prior in parameter space
- Samples shown in data space
- Receive data points (blue circles in data space)
- Compute likelihood
- Posterior is prior (or prev. posterior) times likelihood
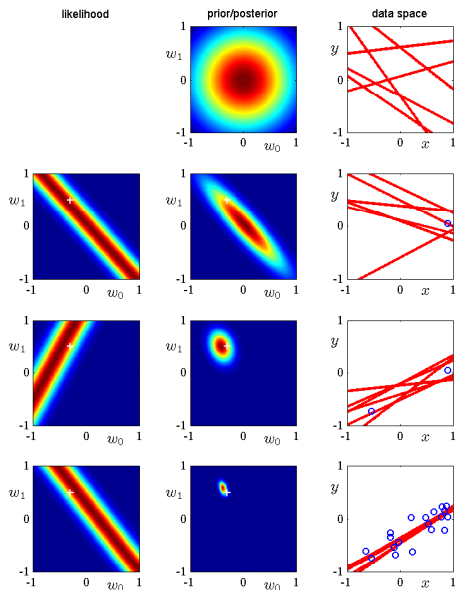
# Bayesian Linear Regression - Intuition



- Simple example $x, t \in \mathbb{R}$, $y(x, \boldsymbol{w}) = w_0 + w_1 x$
- Start with Gaussian prior in parameter space
- Samples shown in data space
- Receive data points (blue circles in data space)
- Compute likelihood
- Posterior is prior (or prev. posterior) times likelihood

# Bayesian Linear Regression - Intuition



- Simple example $x, t \in \mathbb{R}$, $y(x, \boldsymbol{w}) = w_0 + w_1 x$
- Start with Gaussian prior in parameter space
- Samples shown in data space
- Receive data points (blue circles in data space)
- Compute likelihood
- Posterior is prior (or prev. posterior) times likelihood

# Predictive Distribution

- Single estimate of $w$ (ML or MAP) doesn't tell whole story
- We have a distribution over $w$, and can use it to make predictions
- Given a new value for $x$, we can compute a *distribution* over $t$:

$$p(t|\boldsymbol{t}, \alpha, \beta) = \int p(t, \boldsymbol{w}|\boldsymbol{t}, \alpha, \beta)d\boldsymbol{w}$$

$$p(t|\boldsymbol{t}, \alpha, \beta) = \int \underbrace{p(t|\boldsymbol{w}, \beta)}_{predict} \underbrace{p(\boldsymbol{w}|\boldsymbol{t}, \alpha, \beta)}_{probability} \underbrace{d\boldsymbol{w}}_{sum}$$

  - i.e. For each value of $w$, let it make a prediction, multiply by its probability, sum over all $w$
  - For arbitrary models as the distributions, this integral may not be computationally tractable

# Predictive Distribution

- Single estimate of $w$ (ML or MAP) doesn't tell whole story
- We have a distribution over $w$, and can use it to make predictions
- Given a new value for $x$, we can compute a *distribution* over $t$:

$$p(t|\boldsymbol{t}, \alpha, \beta) = \int p(t, \boldsymbol{w}|\boldsymbol{t}, \alpha, \beta) d\boldsymbol{w}$$

$$p(t|\boldsymbol{t}, \alpha, \beta) = \int \underbrace{p(t|\boldsymbol{w}, \beta)}_{predict} \underbrace{p(\boldsymbol{w}|\boldsymbol{t}, \alpha, \beta)}_{probability} \underbrace{d\boldsymbol{w}}_{sum}$$

  - i.e. For each value of $w$, let it make a prediction, multiply by its probability, sum over all $w$
  - For arbitrary models as the distributions, this integral may not be computationally tractable
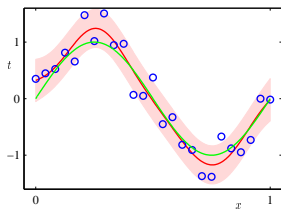
## Predictive Distribution

- Single estimate of $w$ (ML or MAP) doesn't tell whole story
- We have a distribution over $w$, and can use it to make predictions
- Given a new value for $x$, we can compute a *distribution* over $t$:

$$p(t|\boldsymbol{t}, \alpha, \beta) = \int p(t, \boldsymbol{w}|\boldsymbol{t}, \alpha, \beta) d\boldsymbol{w}$$

$$p(t|\boldsymbol{t}, \alpha, \beta) = \int \underbrace{p(t|\boldsymbol{w}, \beta)}_{predict} \underbrace{p(\boldsymbol{w}|\boldsymbol{t}, \alpha, \beta)}_{probability} \underbrace{d\boldsymbol{w}}_{sum}$$

- i.e. For each value of $w$, let it make a prediction, multiply by its probability, sum over all $w$
- For arbitrary models as the distributions, this integral may not be computationally tractable
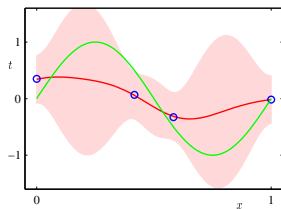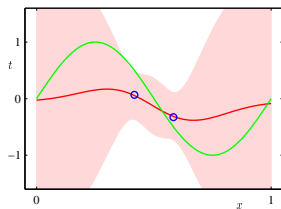
# Predictive Distribution



- With the Gaussians we've used for these distributions, the predicitve distribution will also be Gaussian
  - (math on convolutions of Gaussians)
- Green line is true (unobserved) curve, blue data points, red line is mean, pink one standard deviation
  - Uncertainty small around data points
  - Pink region shrinks with more data

# Bayesian Model Selection

- So what do the Bayesians say about model selection?
    - Model selection is choosing model $\mathcal{M}_i$ e.g. degree of polynomial, type of basis function $\phi$
- Don't select, just integrate

$$p(t|\boldsymbol{x}, \mathcal{D}) = \sum_{i=1}^{L} \underbrace{p(t|\boldsymbol{x}, \mathcal{M}_i, \mathcal{D})}_{predictive\ dist.} p(\mathcal{M}_i|\mathcal{D})$$

- Average together the results of all models
- Could choose most likely model a posteriori $p(\mathcal{M}_i|\mathcal{D})$
    - More efficient, approximation

# Bayesian Model Selection

- So what do the Bayesians say about model selection?
  - Model selection is choosing model $\mathcal{M}_i$ e.g. degree of polynomial, type of basis function $\phi$
- Don't select, just integrate

$$p(t|\boldsymbol{x}, \mathcal{D}) = \sum_{i=1}^{L} \underbrace{p(t|\boldsymbol{x}, \mathcal{M}_i, \mathcal{D})}_{predictive\ dist.} p(\mathcal{M}_i|\mathcal{D})$$

- Average together the results of all models
- Could choose most likely model a posteriori $p(\mathcal{M}_i|\mathcal{D})$
  - More efficient, approximation

# Bayesian Model Selection

- So what do the Bayesians say about model selection?
  - Model selection is choosing model $\mathcal{M}_i$ e.g. degree of polynomial, type of basis function $\phi$
- Don't select, just integrate

$$p(t|\boldsymbol{x}, \mathcal{D}) = \sum_{i=1}^{L} \underbrace{p(t|\boldsymbol{x}, \mathcal{M}_i, \mathcal{D})}_{predictive\ dist.} p(\mathcal{M}_i|\mathcal{D})$$

- Average together the results of all models
- Could choose most likely model a posteriori $p(\mathcal{M}_i|\mathcal{D})$
  - More efficient, approximation

# Bayesian Model Selection

- So what do the Bayesians say about model selection?
    - Model selection is choosing model $\mathcal{M}_i$ e.g. degree of polynomial, type of basis function $\phi$
- Don't select, just integrate

$$p(t|\boldsymbol{x}, \mathcal{D}) = \sum_{i=1}^{L} \underbrace{p(t|\boldsymbol{x}, \mathcal{M}_i, \mathcal{D})}_{\text{predictive dist.}} p(\mathcal{M}_i|\mathcal{D})$$

- Average together the results of all models
- Could choose most likely model a posteriori $p(\mathcal{M}_i|\mathcal{D})$
    - More efficient, approximation

# Bayesian Model Selection

- How do we compute the posterior over models?

$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{D}|\mathcal{M}_i)p(\mathcal{M}_i)$$

- Another likelihood + prior combination
- Likelihood:

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\boldsymbol{w}, \mathcal{M}_i)p(\boldsymbol{w}|\mathcal{M}_i)d\boldsymbol{w}$$

# Bayesian Model Selection

- How do we compute the posterior over models?

$$p(\mathcal{M}_i|\mathcal{D}) \propto p(\mathcal{D}|\mathcal{M}_i)p(\mathcal{M}_i)$$

- Another likelihood + prior combination
- Likelihood:

$$p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\boldsymbol{w}, \mathcal{M}_i)p(\boldsymbol{w}|\mathcal{M}_i)d\boldsymbol{w}$$

# Conclusion

- Readings: Ch. 3.1, 3.1.1-3.1.4, 3.3.1-3.3.2, 3.4
- Linear Models for Regression
    - Linear combination of (non-linear) basis functions
- Fitting parameters of regression model
    - Least squares
    - Maximum likelihood (can be = least squares)
- Controlling over-fitting
    - Regularization
    - Bayesian, use prior (can be = regularization)
- Model selection
    - Cross-validation (use held-out data)
    - Bayesian (use model evidence, likelihood)