

Linear Models for Regression

Greg Mori - CMPT 419/726

Bishop PRML Ch. 3

Outline

Regression

Linear Basis Function Models

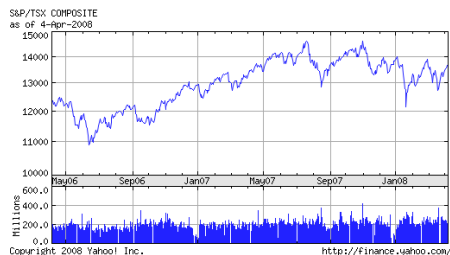
Loss Functions for Regression

Finding Optimal Weights

Regularization

Bayesian Linear Regression

Regression



- Given training set $\{(x_1, t_1), \dots, (x_N, t_N)\}$
- t_i is continuous: regression
- For now, assume $t_i \in \mathbb{R}$, $x_i \in \mathbb{R}^D$
- E.g. t_i is stock price, x_i contains company profit, debt, cash flow, gross sales, number of spam emails sent, ...

Linear Functions

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

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

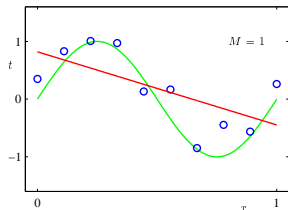
- Linear functions will lead to simple algorithms, so let's see what we can do with them

Linear Regression

- Simplest linear model for regression

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D$$

- Remember, we're learning \mathbf{w}
- Set \mathbf{w} so that $y(\mathbf{x}, \mathbf{w})$ 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

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D$$

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

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

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

- Linear combinations of these **basis functions** also linear in parameters

Linear Basis Function Models

- **Bias** parameter allows fixed offset in data

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

- Think of simple 1-D \mathbf{x} :

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

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

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

Linear Basis Function Models

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

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

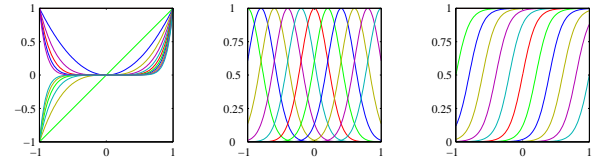
- Polynomial regression is an example of this
- Order M polynomial regression, $\phi_j(x) = ?$
- $\phi_j(x) = x^j$:

$$y(\mathbf{x}, \mathbf{w}) = w_0x^0 + w_1x^1 + \dots + w_Mx^M$$

Basis Functions: Feature Functions

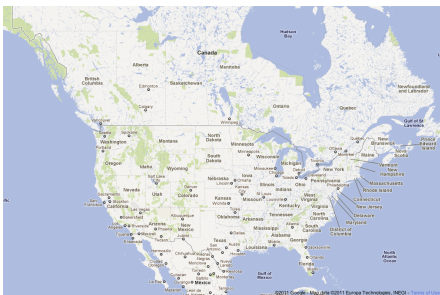
- Often we extract features from x
 - An intuitive way to think of $\phi_j(x)$ is as feature functions
- E.g. Automatic CMPT726 project report grading system
 - x is text of report: In this project we apply the algorithm of Mori [2] to recognizing blue objects. We test this algorithm on pictures of you 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



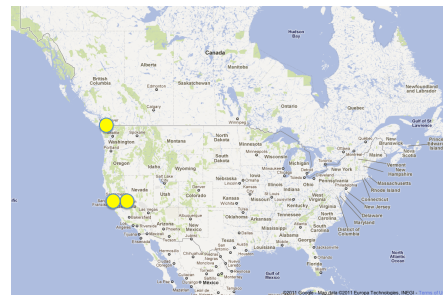
- Polynomial $\phi_j(x) = x^j$
- Gaussians $\phi_j(x) = \exp\{-\frac{(x-\mu_j)^2}{2s^2}\}$
- Sigmoidal $\phi_j(x) = \frac{1}{1+\exp((\mu_j-x)/s)}$

Example - Gaussian Basis Functions: Temperature



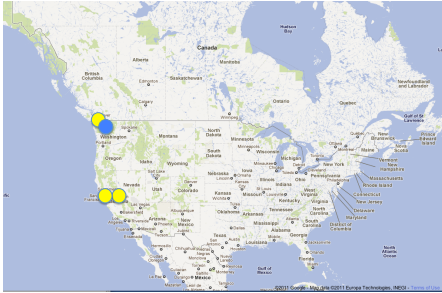
- Use Gaussian basis functions, regression on temperature

Example - Gaussian Basis Functions: Temperature



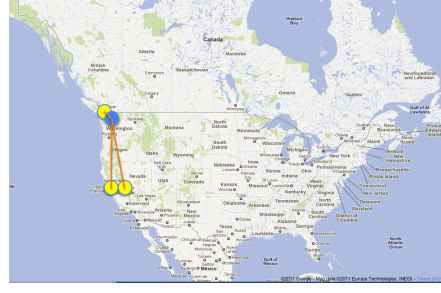
- $\mu_1 = \text{Vancouver}$, $\mu_2 = \text{San Francisco}$, $\mu_3 = \text{Oakland}$

Example - Gaussian Basis Functions: Temperature



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

Example - Gaussian Basis Functions: Temperature



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

Example - Gaussian Basis Functions: 726 Report Grading

- Define:
 - $\mu_1 = \text{Crime and Punishment}$
 - $\mu_2 = \text{Animal Farm}$
 - $\mu_3 = \text{Some paper by Mori}$
- Learn weights:
 - $w_1 = ?$
 - $w_2 = ?$
 - $w_3 = ?$
- Grade a project report x :
 - Measure similarity of x to each μ_i , Gaussian, with weights:

$$y(\mathbf{x}, \mathbf{w}) = w_1 \exp\left\{-\frac{(x-\mu_1)^2}{2s^2}\right\} + w_2 \exp\left\{-\frac{(x-\mu_2)^2}{2s^2}\right\} + w_3 \exp\left\{-\frac{(x-\mu_3)^2}{2s^2}\right\}$$
- The Gaussian basis function models end up similar to template matching

Loss Functions for Regression

- We want to find the “best” set of coefficients \mathbf{w}
- Recall, one way to define “best” was minimizing squared error:

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

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

Gaussian Noise Model for Regression

- We are provided with a training set $\{(x_i, t_i)\}$
- Let's assume t arises from a deterministic function plus Gaussian distributed (with precision β) noise:

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

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

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

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

Finding Optimal Weights

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

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

- In vector form:

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

Maximum Likelihood for Regression

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

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

- The log-likelihood is:

$$\begin{aligned} \ln p(\mathbf{t}|\mathbf{w}, \beta) &= \ln \prod_{n=1}^N \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2} (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2\right) \\ &= \underbrace{\frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)}_{\text{const. wrt } \mathbf{w}} - \beta \underbrace{\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{squared error}} \end{aligned}$$

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

Finding Optimal Weights

- Set gradient to 0:

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

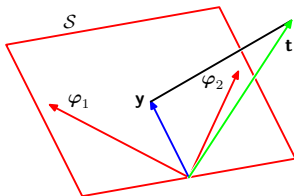
- Maximum likelihood estimate for \mathbf{w} :

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

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

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

Geometry of Least Squares



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

Sequential Learning

- In practice N might be huge, or data might arrive online
- Can use a **gradient descent** method:
 - Start with initial guess for \mathbf{w}
 - Update by taking a step in gradient direction ∇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 ∇E_n for one example
 - Details about step size are important – decrease step size at the end

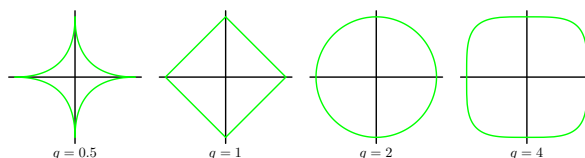
Regularization

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

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

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

Other Regularizers



- Can use different norms for regularizer:

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

- 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:

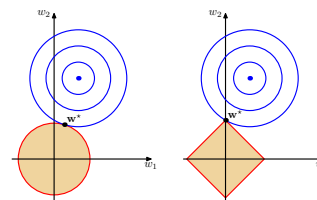
$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Calculus ...

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

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

Ridge Regression vs. Lasso



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

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 \mathbf{w}
- There will turn out to be a connection to regularization

Bayesian Linear Regression

- Start with a prior over parameters \mathbf{w}
 - **Conjugate prior** is a Gaussian:

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

- This simple form will make math easier; can be done for arbitrary Gaussian too
- Data likelihood, Gaussian model as before:

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

Bayesian Linear Regression

- Posterior distribution on \mathbf{w} :

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

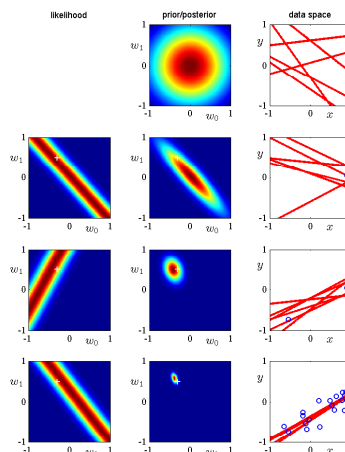
$$= \left[\prod_{n=1}^N \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2}(t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2\right) \right] \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \exp\left(-\frac{\alpha}{2}\mathbf{w}^T \mathbf{w}\right)$$

- Take the log:

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

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

Bayesian Linear Regression - Intuition



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

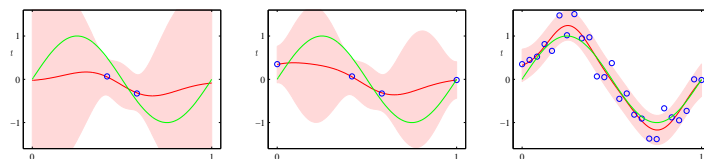
- Single estimate of \mathbf{w} (ML or MAP) doesn't tell whole story
- We have a distribution over \mathbf{w} , and can use it to make predictions
- Given a new value for x , we can compute a *distribution* over t :

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

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

- i.e. For each value of \mathbf{w} , let it make a prediction, multiply by its probability, sum over all \mathbf{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 predictive 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 ϕ
- Don't select, just integrate

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

- 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:

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

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)