#### Lecture 4

- Homework
  - Hw 1 and 2 will be reoped after class for every body. New deadline 4/20
  - Hw 3 and 4 online (Nima is lead)
- Pod-cast lecture on-line
- Final projects
  - Nima will register groups next week. Email/tell Nima.
  - Give proposal in week 5
  - See last years topic on webpage. Choose your own or
- Linear regression 3.0-3.3+ SVD
- Next lectures:
  - I posted a rough plan.
  - It is flexible though so please come with suggestions

#### Projects

- **3-4** person groups
- Deliverables: Poster & Report & main code (plus proposal, midterm slide)
- Topics your own or chose form suggested topics
- Week 3 groups due to TA Nima (if you don't have a group, ask in week 2 and we can help).
- Week 5 proposal due. TAs and Peter can approve.
- Proposal: One page: Title, A large paragraph, data, weblinks, references.
- Something physical
- Week ~7 Midterm slides? Likely presented to a subgroup of class.
- Week 10/11 (likely 5pm 6 June Jacobs Hall lobby) final poster session?
- Report due Saturday 16 June.

#### Mark's Probability and Data homework

## Mark's Probability and Data homework

### Linear regression: Linear Basis Function Models (1)

Generally

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

- where  $\phi_j(x)$  are known as *basis functions*.
- Typically,  $\phi_0(x) = 1$ , so that  $w_0$  acts as a bias.
- Simplest case is linear basis functions:  $\phi_d(x) = x_d$ . -1



http://playground.tensorflow.org/

#### Maximum Likelihood and Least Squares (3)

Computing the gradient and setting it to zero yields

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t} | \mathbf{w}, \beta) = \beta \sum_{n=1}^{N} \left\{ t_n - \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n) \right\} \boldsymbol{\phi}(\mathbf{x}_n)^{\mathrm{T}} = \mathbf{0}.$$
Solving for w,
$$\mathbf{w}_{\mathrm{ML}} = \left( \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi} \right)^{-1} \mathbf{\Phi}^{\mathrm{T}} \mathbf{t}$$

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

Least mean squares: An alternative approach for big datasets



This is **"on-line" learning**. It is efficient if the dataset is redundant and simple to implement.

- It is called stochastic gradient descent if the training cases are picked randomly.
- Care must be taken with the learning rate to prevent divergent oscillations. Rate must decrease with tau to get a good fit.

**Regularized least squares** 

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{ y(\mathbf{x}_n, \mathbf{w}) - t_n \}^2 + \frac{\lambda}{2} \| \mathbf{w} \|^2$$

The squared weights penalty is mathematically compatible with the squared error function, giving a closed form for the optimal weights:

$$\mathbf{w}^* = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$
identity matrix

#### A picture of the effect of the regularizer



- The overall cost function is the sum of two parabolic bowls.
- The sum is also a parabolic bowl.
- The combined minimum lies on the line between the minimum of the squared error and the origin.
- The L2 regularizer just **shrinks** the weights.

#### Other regularizers

- We do not need to use the squared error, provided we are willing to do more computation.
- Other powers of the weights can be used.



The lasso: penalizing the absolute values of the weights

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{ y(\mathbf{x}_n, \mathbf{w}) - t_n \}^2 + \lambda \sum_i |\mathbf{w}_i|$$

- Finding the minimum requires quadratic programming but its still unique because the cost function is convex (a bowl plus an inverted pyramid)
- As lambda increases, many weights go to exactly zero.
  - This is great for interpretation, and it is also prevents overfitting.

# Geometrical view of the lasso compared with a penalty on the squared weights





optimum

Minimizing the absolute error

$$\min_{over \mathbf{w}} \sum_{n} |t_n - \mathbf{w}^T \mathbf{x}_n|$$

- This minimization involves solving a linear programming problem.
- It corresponds to maximum likelihood estimation if the output noise is modeled by a Laplacian instead of a Gaussian.

$$p(t_n | y_n) = a e^{-a |t_n - y_n|} -\log p(t_n | y_n) = -a |t_n - y_n| + const$$

## The bias-variance trade-off (a figment of the frequentists lack of imagination?)

- Imagine a training set drawn at random from a whole set of training sets.
- The squared loss can be decomposed into a
  - Bias = systematic error in the model's estimates
  - Variance = noise in the estimates cause by sampling noise in the training set.
- There is also additional loss due to noisy target values.
  - We eliminate this extra, irreducible loss from the math by using the average target values (i.e. the unknown, noise-free values)



#### The bias-variance decomposition



Regularization parameter affects the bias and variance terms



### An example of the bias-variance trade-off



#### Beating the bias-variance trade-off

- Reduce the variance term by averaging lots of models trained on different datasets.
  - Seems silly. For lots of different datasets it is better to combine them into one big training set.
    - With more training data there will be much less variance.
- Weird idea: We can create different datasets by bootstrap sampling of our single training dataset.
  - This is called "bagging" and it works surprisingly well.
- If we have enough computation its better doing it **Bayesian**:
  - Combine the predictions of many models using the posterior probability of each parameter vector as the combination weight.

#### **Bayesian** Linear Regression (1)

Define a conjugate prior over w

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

•Combining this with the likelihood function and using results for marginal and conditional Gaussian distributions, gives the posterior

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \qquad \mathbf{m}_N = \mathbf{S}_N \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \mathbf{\Phi}^{\mathrm{T}} \mathbf{t} \right) \\ \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi}.$$

• A common simpler prior

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

• Which gives

$$\mathbf{m}_N = eta \mathbf{S}_N \mathbf{\Phi}^{\mathrm{T}} \mathbf{t} \ \mathbf{S}_N^{-1} = lpha \mathbf{I} + eta \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi}.$$

#### From lecture 3:

Bayes for linear model y = Ax + n  $n \sim N(0, C_n)$   $y \sim N(Ax, C_n)$  prior:  $x \sim N(0, C_x)$  $x_{p} = C_{p}A^{T}C_{n}^{-1}y$   $= C_{2}(y - A_{x})^{T}C_{n}^{-1}(y - A_{x}) = C_{2}(y - A_{x})^{T}C_{n}^{-1}(y - A_{x}) = C_{2}(x - A_{x})^{T}C_{n}^{-1}$  $p(\mathbf{x}|\mathbf{y}) \sim p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) \sim N(\mathbf{x}_{\mathbf{p}}, \mathbf{C}_{n})$  $n p^{-\frac{1}{2}(x-x_p)} (x-x_p) =$  $C_{\varphi}^{-1} = A^{\dagger}C_{\varphi}^{-1}A + C_{\chi}^{\dagger}$ xp= CpA'Ch'q

#### Interpretation of solution

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

Draw it

Sequential, conjugate prior

$$p(\mathbf{x}|\mathbf{y}) \sim p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) \sim N(\mathbf{A}\mathbf{x}, \mathbf{C}_n) \ N(\mathbf{0}, \mathbf{C}_{\mathbf{x}}) \sim N(\mathbf{x}_p, \mathbf{C}_p)$$
  
Covariance  $\mathbf{C}_p^{-1} = \mathbf{A}^T \mathbf{C}_n^{-1} \mathbf{A} + \mathbf{C}_x^{-1}$ 



#### **Predictive Distribution**

Predict t for new values of x by integrating over w (Giving the marginal distribution of t):

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

$$\uparrow \qquad \uparrow \qquad = \mathcal{N}(t|\mathbf{m}_N^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$
training data
precision of prior
precision of output noise

$$\mathbf{m}_N = eta \mathbf{S}_N \mathbf{\Phi}^{\mathrm{T}} \mathbf{t} \ \mathbf{S}_N^{-1} = lpha \mathbf{I} + eta \mathbf{\Phi}^{\mathrm{T}} \mathbf{\Phi}.$$

• where

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

$$p(t|\mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) \, \mathrm{d}\mathbf{w}$$
$$= \mathcal{N}(t|\mathbf{m}_N^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x}))$$

• Just use ML solution

• Prior predictive

# Predictive distribution for noisy sinusoidal data modeled by linear combining 9 radial basis functions.



A way to see the covariance of predictions for different values of x

We sample models at random from the posterior and *show the mean* of each model's predictions



# Equivalent Kernel BISHOP 3.3.3

The predictive mean can be written

$$y(\mathbf{x}, \mathbf{m}_{N}) = \mathbf{m}_{N}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}) = \beta \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \mathbf{S}_{N} \boldsymbol{\Phi}^{\mathrm{T}} \mathbf{t}$$
$$= \sum_{n=1}^{N} \beta \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \mathbf{S}_{N} \boldsymbol{\phi}(\mathbf{x}_{n}) t_{n}$$
$$= \sum_{n=1}^{N} k(\mathbf{x}, \mathbf{x}_{n}) t_{n}.$$
Equivalent kernel or smoother matrix.

This is a weighted sum of the training data target values,  $t_n$ .

#### Equivalent Kernel (2)



Weight of  $t_n$  depends on distance between x and  $x_n$ ; nearby  $x_n$  carry more weight.

#### Equivalent Kernel (4)

• The kernel as a covariance function: consider

$$\begin{aligned} \operatorname{cov}[y(\mathbf{x}), y(\mathbf{x}')] &= \operatorname{cov}[\boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \mathbf{w}, \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}')] \\ &= \boldsymbol{\phi}(\mathbf{x})^{\mathrm{T}} \mathbf{S}_{N} \boldsymbol{\phi}(\mathbf{x}') = \beta^{-1} k(\mathbf{x}, \mathbf{x}') \end{aligned}$$

- We can avoid the use of basis functions and define the kernel function directly, leading to *Gaussian Processes* (Chapter 6).
- No need to determine weights.

• Like all kernel functions, the equivalent kernel can be expressed as an inner product:

$$egin{aligned} k(\mathbf{x},\mathbf{z}) &= oldsymbol{\psi}(\mathbf{x})^{\mathrm{T}}oldsymbol{\psi}(\mathbf{z}) \ oldsymbol{\psi}(\mathbf{x}) &= eta^{1/2} \mathbf{S}_N^{1/2} oldsymbol{\phi}(\mathbf{x}) \end{aligned}$$

#### SVD

y = Xw

$$\boldsymbol{X} = \boldsymbol{U}\boldsymbol{S}\boldsymbol{V}^T = [\boldsymbol{u}_1, \dots, \boldsymbol{u}_N]\boldsymbol{S} [\boldsymbol{v}_1, \dots, \boldsymbol{v}_M]^T$$

 $\begin{array}{l} \boldsymbol{X}\boldsymbol{X}^T = \\ \boldsymbol{X}^T\boldsymbol{X} = \end{array}$ 

$$w_{LS} = (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X} \boldsymbol{y} =$$

$$w_{Reg} = (\lambda I + X^T X)^{-1} X y =$$