

# Regression and Prediction

Class 15. 23 Oct 2012

Instructor: Bhiksha Raj

- The derivative of a scalar function w.r.t. a vector is a vector
- The derivative w.r.t. a matrix is a matrix

$$df(\mathbf{x}) = \begin{bmatrix} \frac{df}{dx_1} & \frac{df}{dx_2} & \dots & \frac{df}{dx_p} \end{bmatrix}$$

## Matrix Identities

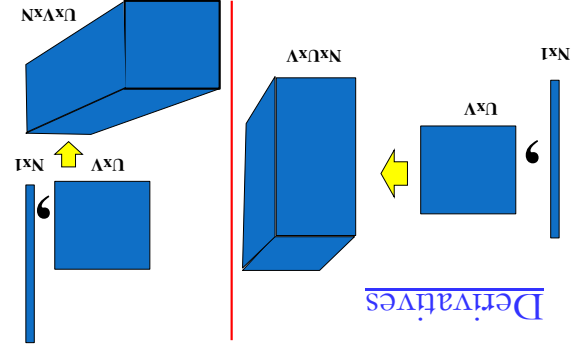
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_p \end{bmatrix}$$

## Matrix Identities

$$df(\mathbf{x}) = \begin{bmatrix} \frac{df}{dx_{p1}} & \frac{df}{dx_{p2}} & \dots & \frac{df}{dx_{po}} \\ \frac{df}{dx_{21}} & \frac{df}{dx_{22}} & \dots & \frac{df}{dx_{2o}} \\ \frac{df}{dx_{11}} & \frac{df}{dx_{12}} & \dots & \frac{df}{dx_{1o}} \end{bmatrix}$$

- The derivative of a scalar function w.r.t. a vector is a vector
- The derivative w.r.t. a matrix is a matrix

## Derivatives



- In general: Differentiating an MXN function by a UVL argument results in an MxNxUxV tensor derivative

## Matrix derivative identities

$$d(\mathbf{X}\mathbf{a}) = \mathbf{X}d\mathbf{a} \quad d(\mathbf{a}^T\mathbf{X}) = \mathbf{X}^T d\mathbf{a}$$

$$d(\mathbf{A}\mathbf{X}) = (d\mathbf{A})\mathbf{X} ; \quad d(\mathbf{X}\mathbf{A}) = \mathbf{X}(d\mathbf{A})$$

**X** is a matrix, **a** is a vector  
Solution may also be  $\mathbf{X}^T$   
**A** is a matrix

$$d(\mathbf{a}^T\mathbf{X}\mathbf{a}) = \mathbf{a}^T d\mathbf{X} + \mathbf{X}^T d\mathbf{a}$$

$$d(\mathbf{A}\mathbf{X}\mathbf{A}^T) = d\mathbf{A}\mathbf{X}\mathbf{A}^T + \mathbf{A}d\mathbf{X}\mathbf{A}^T + \mathbf{A}^T d\mathbf{A}\mathbf{X}$$

- Some basic linear and quadratic identities

- The derivative of a vector function w.r.t. a vector is a matrix
- Note transposition of order

$$\mathbf{F}(\mathbf{x}) = \begin{bmatrix} F_1 \\ F_2 \\ \dots \\ F_N \end{bmatrix} = \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_p \end{bmatrix}$$

$$d\mathbf{F}(\mathbf{x}) = \begin{bmatrix} \frac{dF_1}{dx_1} & \frac{dF_1}{dx_2} & \dots & \frac{dF_1}{dx_p} \\ \frac{dF_2}{dx_1} & \frac{dF_2}{dx_2} & \dots & \frac{dF_2}{dx_p} \\ \dots & \dots & \dots & \dots \\ \frac{dF_N}{dx_1} & \frac{dF_N}{dx_2} & \dots & \frac{dF_N}{dx_p} \end{bmatrix}$$

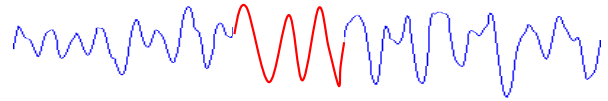
## Matrix Identities

### ■ Generally a tool to *predict* variables

- Linear regression, Simple regression, Ordinary least squares, Polynomial regression, General linear model, Generalized linear model, Discrete choice, Logistic regression, Multinomial logit, Mixed logit, Probit, Multinomial probit, ....
- Wikipedia
- Expressed in many forms
- Analyzing relationship between variables

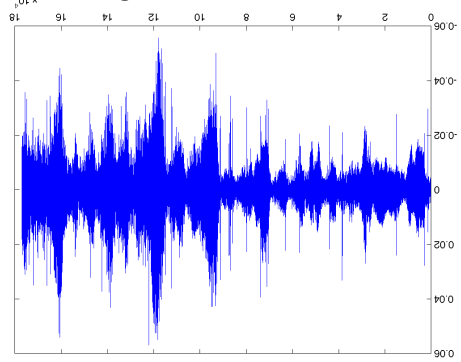
## What is a regression

- "Extend" the curve on the left to "predict" the values in the "blank" region
- *Forward prediction*
- Extend the blue curve on the right leftwards to predict the blank region
- *Backward prediction*
- How?
- Regression analysis..



## Interpolation..

### ■ Can you spot the glitches?

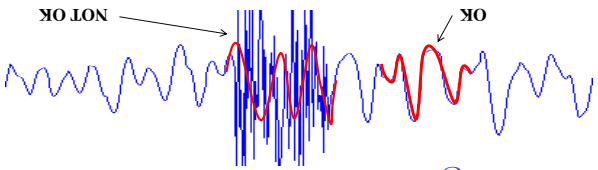


## A Common Problem

## Regressions for prediction

- $y = f(x; \Theta) + e$
- Different possibilities
  - $y$  is a scalar
  - $y$  is real
  - $y$  is categorical (classification)
  - $y$  is a vector
  - $x$  is a vector
  - $x$  is a set of real valued variables
  - $x$  is a set of categorical variables
  - $x$  is a combination of the two
- $f(\cdot)$  is a linear or affine function
- $f(\cdot)$  is a non-linear function
- $f(\cdot)$  is a *time-series* model

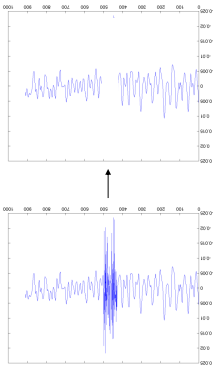
- Regression-based reconstruction can be done anywhere
- Reconstructed value will not match actual value
- Large error of reconstruction identifies glitches



## Detecting the Glitch

## How to fix this problem?

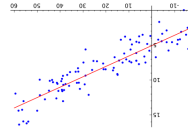
- "Glitches" in audio
  - Must be detected
  - How?
- Then what?
  - Delete the glitch
  - Results in a "hole"
  - Fill in the hole
  - How?
- Glitches must be "fixed"



- Given training data: several  $\mathbf{x}, \mathbf{y}$
- Can define a "divergence":  $D(\mathbf{y}, \hat{\mathbf{y}})$
- Measures how much what differs from  $\mathbf{y}$
- Ideally, if the model is accurate this should be small
- Estimate  $\mathbf{A}, \mathbf{b}$  to minimize  $D(\mathbf{y}, \hat{\mathbf{y}})$

Assuming no error

$$\hat{\mathbf{y}} = \mathbf{A}^T \mathbf{X}$$

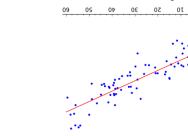
$$\mathbf{y} = \mathbf{A}^T \mathbf{X} + \mathbf{e}$$


## Learning the parameters

- Given a "training" set of  $\{\mathbf{x}, \mathbf{y}\}$  values: estimate  $\mathbf{A}$  and  $\mathbf{b}$
- $\mathbf{y}_1 = \mathbf{A}\mathbf{x}_1 + \mathbf{b} + \mathbf{e}_1$
- $\mathbf{y}_2 = \mathbf{A}\mathbf{x}_2 + \mathbf{b} + \mathbf{e}_2$
- $\mathbf{y}_3 = \mathbf{A}\mathbf{x}_3 + \mathbf{b} + \mathbf{e}_3$
- ...
- If  $\mathbf{A}$  and  $\mathbf{b}$  are well estimated, prediction error will be small

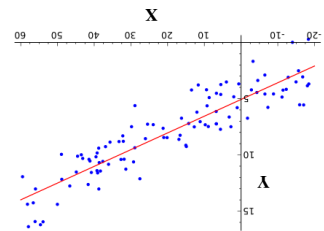
$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} + \mathbf{e}$$

$\mathbf{e}$  = prediction error



## Linear Regressions

- Assumption: relationship between variables is linear
- A linear trend may be found relating  $\mathbf{x}$  and  $\mathbf{y}$
- $\mathbf{y}$  = dependent variable
- $\mathbf{x}$  = explanatory variable
- Given  $\mathbf{x}, \mathbf{y}$  can be predicted as an affine function of  $\mathbf{x}$



## A linear regression

- Define the divergence as the sum of the squared error in predicting  $\mathbf{y}$

$$\mathbf{E} = (\mathbf{y} - \mathbf{A}^T \mathbf{X})(\mathbf{y} - \mathbf{A}^T \mathbf{X})^T = \|\mathbf{y} - \mathbf{A}^T \mathbf{X}\|^2$$

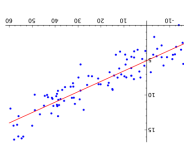
$$D(\mathbf{y}, \hat{\mathbf{y}}) = \mathbf{E} = e_1^2 + e_2^2 + e_3^2 + \dots$$

$$= (y_1 - \mathbf{a}^T \mathbf{b})^2 + (y_1 - \mathbf{a}^T \mathbf{b})^2 + (y_1 - \mathbf{a}^T \mathbf{b})^2 + \dots$$

$$\mathbf{y} = \mathbf{A}^T \mathbf{X} + \mathbf{e}$$

$$y_1 = \mathbf{a}^T \mathbf{x}_1 + \mathbf{b} + e_1$$

$$y_2 = \mathbf{a}^T \mathbf{x}_2 + \mathbf{b} + e_2$$

$$y_3 = \mathbf{a}^T \mathbf{x}_3 + \mathbf{b} + e_3$$


## The prediction error as divergence

$$\mathbf{y} = \mathbf{A}^T \mathbf{X} + \mathbf{e}$$

- Rewrite

$$\mathbf{y} = [y_1 \ y_2 \ y_3 \dots]$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \dots \\ 1 & 1 & 1 & \dots \end{bmatrix}$$

$$\mathbf{A} = \begin{bmatrix} \mathbf{a} \\ b \end{bmatrix}$$

$$\mathbf{e} = [e_1 \ e_2 \ e_3 \dots]$$

- Define:

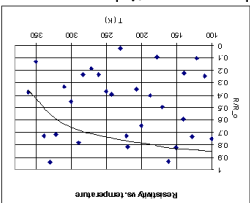
$$y_1 = \mathbf{a}^T \mathbf{x}_1 + \mathbf{b} + e_1$$

$$y_2 = \mathbf{a}^T \mathbf{x}_2 + \mathbf{b} + e_2$$

$$y_3 = \mathbf{a}^T \mathbf{x}_3 + \mathbf{b} + e_3$$

## Linear Regression to a scalar

Check this shit out (Fig. 1). That's bonafide, 100%-real data, my friends. I took it myself over the course of two weeks. And this was not a leisurely two weeks, either; I busted my ass day and night in order to provide you with nothing but the best data possible. Now, let's look a bit more closely at this data, remembering that it is absolutely first-rate. Do you see the exponential dependence? I sure don't. I see a bunch of crap. Christ, this was such a waste of my time. Banking on my hopes that whoever grades this will just look at the pictures, I drew an exponential through my noise. I believe the apparent legitimacy is enhanced by the fact that I used a complicated computer program to make the fit. I understand this is the same process by which the top quark was discovered.



## An imaginary regression.

### Multiple Regression

$$Y = [y_1 \ y_2 \ y_3 \dots] \quad E = [e_1 \ e_2 \ e_3 \dots]$$

$$X = \begin{bmatrix} x_1 & x_2 & x_3 & \dots \\ 1 & 1 & 1 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix} \quad A = \begin{bmatrix} a \\ b \end{bmatrix}$$

[Dx] vector of ones

$$Y = A^T X + E$$

$$DIV = \sum_i \|y_i - A^T x_i - b\|_2^2 = trace(Y - A^T X)(Y - A^T X)^T$$

$$dDIV = (2A^T X X^T - 2Y X^T) dA = 0$$

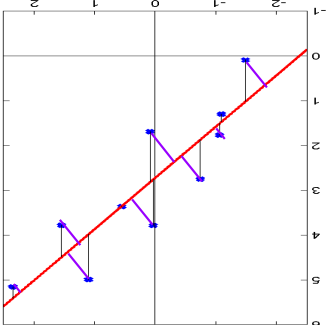
$$A^T = Y X^T (X X^T)^{-1} \quad A = (X X^T)^{-1} X^T Y$$

$$A = A^T X + E$$

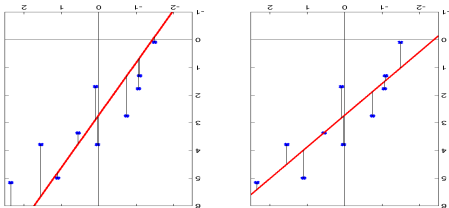
Differentiating and equating to 0

### An Aside

What happens if we minimize the perpendicular instead?



### Prediction error as divergence



$y = a^T x + e$

- Find the "slope"  $a$  such that the total squared length of the error lines is minimized
- $e$  = prediction error

### Solving a linear regression

Minimize squared error

$$E = \|y - X^T A\|_2^2 = (y - A^T X)(y - A^T X)^T$$

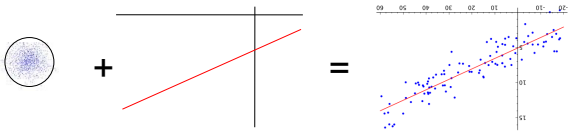
$$= y y^T + A^T X X^T A - 2y X^T A$$

Differentiating w.r.t  $A$  and equating to 0

$$dE = (2A^T X X^T - 2y X^T) dA = 0$$

$$A^T = y X^T (X X^T)^{-1} \quad A = (X X^T)^{-1} X^T y$$

### A Different Perspective



$y$  is a noisy reading of  $A^T x$

$$y = A^T x + e$$

Error  $e$  is Gaussian

$$e \sim N(0, \sigma^2 I)$$

Estimate  $A$  from  $Y = [y_1 \ y_2 \dots y_N]$   $X = [x_1 \ x_2 \dots x_N]$

Fundamentally no different from  $N$  separate single regressions

But we can use the relationship between  $Y$ 's to our benefit

Also called *multiple regression*

Equivalent of saying:

$$y_1 = A^T x_1 + b + e_1$$

$$y_2 = A^T x_2 + b + e_2$$

$$y_3 = A^T x_3 + b + e_3$$

$y_j$  is a vector

$y_j$  =  $j^{\text{th}}$  component of vector  $y_j$

$a$  =  $j^{\text{th}}$  column of  $A$

$b_j$  =  $i^{\text{th}}$  component of  $b$

$$y_{11} = a_1^T x_1 + b_1 + e_{11}$$

$$y_{12} = a_2^T x_2 + b_2 + e_{12}$$

$$y_{13} = a_3^T x_3 + b_3 + e_{13}$$

## The Likelihood of the data

$$y = A^T x + e \quad e \sim N(0, \sigma^2 \mathbf{I})$$

- Probability of observing a specific  $y$ , given  $x$ , for a particular matrix  $A$

$$P(y | x; A) = N(A^T x, \sigma^2 \mathbf{I})$$

- Probability of the collection:  $\mathbf{y} = [y_1, y_2, \dots, y_N]$ ,  $\mathbf{X} = [x_1, x_2, \dots, x_N]$

$$P(\mathbf{y} | \mathbf{X}; A) = \prod_{i=1}^N N(A^T x_i, \sigma^2 \mathbf{I})$$

- Assuming IID for convenience (not necessary)

## The Likelihood of the data

$$y = A^T x + e \quad e \sim N(0, \sigma^2 \mathbf{I})$$

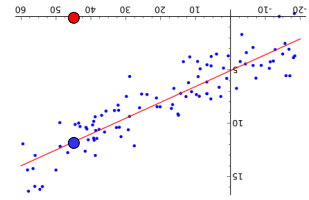
- Probability of observing a specific  $y$ , given  $x$ , for a particular matrix  $A$

$$P(y | x; A) = N(A^T x, \sigma^2 \mathbf{I})$$

- Probability of the collection:  $\mathbf{y} = [y_1, y_2, \dots, y_N]$ ,  $\mathbf{X} = [x_1, x_2, \dots, x_N]$

$$P(\mathbf{y} | \mathbf{X}; A) = \prod_{i=1}^N N(A^T x_i, \sigma^2 \mathbf{I})$$

- Assuming IID for convenience (not necessary)



## Predicting an output

- From a collection of training data, have learned  $A$
- Given  $x$  for a new instance, but not  $y$ , what is  $y$ ?
- Simple solution:

$$\hat{y} = A^T x$$

## A Maximum Likelihood Estimate

$$y = A^T x + e \quad e \sim N(0, \sigma^2 \mathbf{I}) \quad \mathbf{Y} = [y_1, y_2, \dots, y_N] \quad \mathbf{X} = [x_1, x_2, \dots, x_N]$$

$$P(\mathbf{y} | \mathbf{X}) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\sigma^2}{2} \|A^T x_i\|^2\right)$$

$$\log P(\mathbf{y} | \mathbf{X}; A) = C - \sum_{i=1}^N \frac{1}{2\sigma^2} \|y_i - A^T x_i\|^2$$

$$= C - \frac{1}{2\sigma^2} \text{trace}(\mathbf{Y} - A^T \mathbf{X})(\mathbf{Y} - A^T \mathbf{X})^T$$

- Maximizing the log probability is identical to

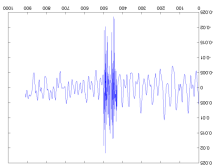
- minimizing the trace

- Identical to the least squares solution

$$A^T = \mathbf{YX}^T (\mathbf{X}\mathbf{X}^T)^{-1} = \mathbf{Y} \text{pinv}(\mathbf{X}) \quad A = (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X}\mathbf{Y}^T$$

## Applying it to our problem

- Prediction by regression

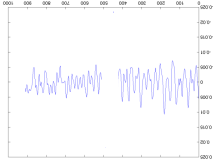


$$x_t = a_1 x_{t-1} + a_2 x_{t-2} + \dots + a_k x_{t-k} + e_t$$

- Forward regression

- Backward regression

$$x_t = b_1 x_{t+1} + b_2 x_{t+2} + \dots + b_k x_{t+k} + e_t$$

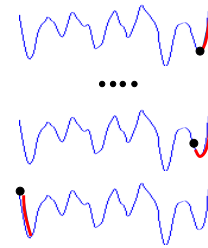


## Applying it to our problem

- Forward prediction

$$\mathbf{x} = \mathbf{a}_t^T \mathbf{X} + e$$

$$\begin{bmatrix} x_t \\ x_{t-1} \\ \dots \\ x_{t-k+1} \end{bmatrix} = \mathbf{a}_t^T \begin{bmatrix} x_{t-1} & x_{t-2} & \dots & x_{t-k} \\ x_{t-2} & x_{t-3} & \dots & x_{t-k-1} \\ \dots & \dots & \dots & \dots \\ x_{t-k-1} & x_{t-k-2} & \dots & x_{t-k-1} \end{bmatrix} + \begin{bmatrix} e_t \\ e_{t-1} \\ \dots \\ e_{t-k+1} \end{bmatrix}$$

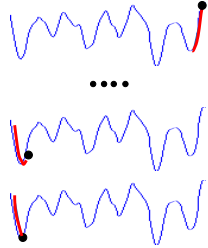


- Backward prediction

$$\underline{\mathbf{x}} = \text{pinv}(\underline{\mathbf{X}}) \mathbf{b}_t$$

$$\underline{\mathbf{x}} = \mathbf{b}_t^T \underline{\mathbf{X}} + e$$

$$\begin{bmatrix} x_{t-k-1} \\ x_{t-k-2} \\ \dots \\ x_{t-k} \\ x_{t-1} \\ x_{t-2} \\ \dots \\ x_{t-1} \\ x_2 \\ \dots \\ e_{t-k-2} \\ e_{t-k-1} \end{bmatrix} = \mathbf{b}_t^T \begin{bmatrix} x_{t-1} & x_{t-2} & \dots & x_{t-k} \\ x_{t-2} & x_{t-3} & \dots & x_{t-k-1} \\ \dots & \dots & \dots & \dots \\ x_{t-k-1} & x_{t-k-2} & \dots & x_{t-k-1} \end{bmatrix} + \begin{bmatrix} e_t \\ e_{t-1} \\ \dots \\ e_{t-k+1} \end{bmatrix}$$



$\hat{y} = \mathbf{Y} \mathbf{X}^T \hat{\mathbf{x}} = \sum_i \mathbf{x}_i^T \hat{\mathbf{x}} y_i$   
 ■ Weighted combination of inputs

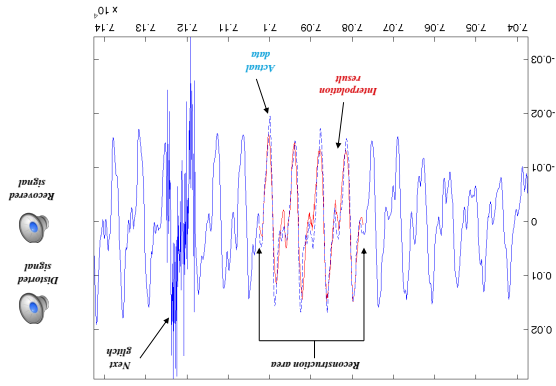
- Let  $\hat{\mathbf{x}} = \mathbf{C}^{-1} \mathbf{x}$
- Normalizing and rotating space
- The rotation is irrelevant

$\mathbf{C} = \mathbf{X} \mathbf{X}^T$

■ What are we doing exactly?

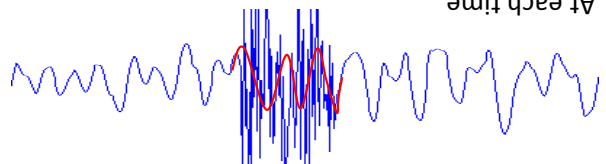
$\mathbf{A} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y}^T$       $\hat{y} = \mathbf{A}^T \mathbf{x} = \mathbf{Y} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{x}$

Predicting a value



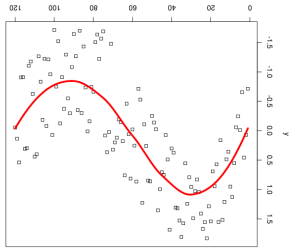
Reconstruction zoom in

- At each time
  - Learn a "forward" predictor  $a_1$
  - Compute error:  $err_t = |x_t^1 - x_{t-1}^{est}|^2$
  - Learn a "backward" predictor and compute backward error
  - $ber_t$
  - Compute average prediction error over window, threshold



Finding the burst

- Multiple solutions
- How do we model these?



Relationships are not always linear

- Note the structure
- Can also be done in batch mode!

$\mathbf{a}^{t+1} = \mathbf{a}^t + \eta (\mathbf{y}_t - \hat{\mathbf{y}}_t) \mathbf{x}_t$       $\hat{\mathbf{y}}_t = (\mathbf{a}^t)^T \mathbf{x}_t$

Scalar prediction version

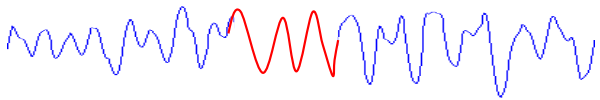
- The Widrow Hoff rule
- As data comes in?
- Can we learn A incrementally instead?

$\mathbf{A} = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y}^T$

Requires knowledge of all (x,y) pairs

Incrementally learning the regression

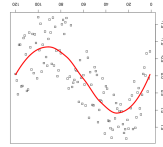
- Learn "forward" predictor at left edge of "hole"
  - For each missing sample
  - At each time, predict next sample  $x_{t+1}^{est} = \sum_i a_{i,t} x_{t-k}^i$
  - Use estimated samples if real samples are not available
- Learn "backward" predictor at left edge of "hole"
  - For each missing sample
  - At each time, predict next sample  $x_{t+1}^{est} = \sum_i b_{i,t} x_{t+k}^i$
  - Use estimated samples if real samples are not available
- Average forward and backward predictions



Filling the hole

- But first.. MAP estimators..
  - The "kernel" is the kernel of a parzen window
  - Note – an estimator of  $y$ , not parameters of regression
- Actually a non-parametric MAP estimator of  $y$

$$\hat{y} = \frac{\sum_i K^h(\mathbf{x} - \mathbf{x}_i) y_i}{\sum_i K^h(\mathbf{x} - \mathbf{x}_i)}$$

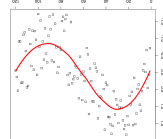


### Kernel Regression

$$\mathbf{y} = \sum_i d(\mathbf{x}, \mathbf{x}_i) \mathbf{y}_i + \mathbf{e}$$

- For any  $\mathbf{x}$
- How about doing this locally?

$$\hat{\mathbf{y}} = \mathbf{Y} \hat{\mathbf{X}}^T \mathbf{x} = \sum_i \mathbf{x}_i^T \hat{\mathbf{y}}_i \mathbf{x}_i \mathbf{y}_i + \mathbf{e}$$



- Must apply everywhere
- Regression is usually trained over the entire data

### Being non-committal: Local Regression

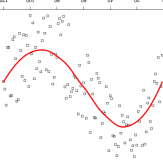
$$\mathbf{A} = (\Phi(\mathbf{X})\Phi(\mathbf{X})^T)^{-1}\Phi(\mathbf{X})\mathbf{Y}^T$$

- $\mathbf{Y} = \mathbf{A}\Phi(\mathbf{X}) + \mathbf{e}$
- Replace  $\mathbf{X}$  with  $\Phi(\mathbf{X})$  in earlier equations for solution

$$\mathbf{X} \leftarrow \Phi(\mathbf{X}) = [\phi(\mathbf{x}_1) \ \phi(\mathbf{x}_2) \ \dots \ \phi(\mathbf{x}_N)]$$

$$\mathbf{x} \leftarrow \phi(\mathbf{x}) = [\phi_1(\mathbf{x}) \ \phi_2(\mathbf{x}) \ \dots \ \phi_N(\mathbf{x})]$$

- $y = \phi(\mathbf{x}) + e$



### Non-linear regression

- MAP is simpler to visualize
- ML (Maximum Likelihood): Find that value of  $y$  for which the statistical best guess of  $X$  would have been the observed  $X$ 

$$y = \text{argmax}_y P(\mathbf{x}|y)$$
- MAP (Maximum A Posteriori): Find a "best guess" for  $y$  (in a statistical sense), given that we know  $\mathbf{x}$ 

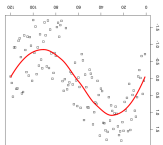
$$y = \text{argmax}_y P(y|\mathbf{x})$$

### Map Estimators

- No closed form solution
  - But can be highly accurate
- But what is  $d(\mathbf{x}, \mathbf{x}_i)$ ??
 
$$e(\mathbf{x}) = \|\mathbf{y} - \sum_i d(\mathbf{x}, \mathbf{x}_i) \mathbf{y}_i\|_2^2$$

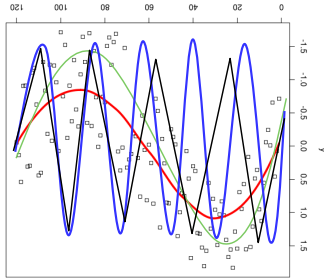
$$\hat{\mathbf{y}} = \sum_i d(\mathbf{x}, \mathbf{x}_i) \mathbf{y}_i$$

- The resulting regression is dependent on  $\mathbf{x}_i$



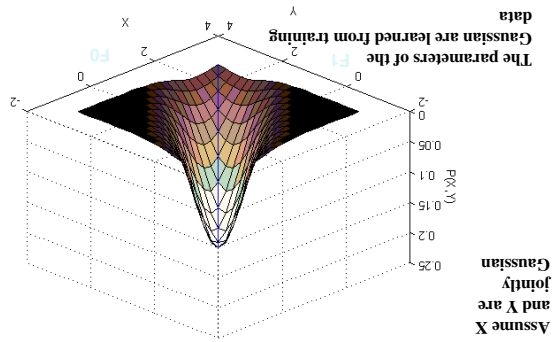
### Local Regression

- Finding the optimal combination of various function
  - Remind you of something?



### What we are doing

## MAP estimation: Gaussian PDF



43

11755/18797

23 Oct 2012

## Learning the parameters of the Gaussian

$$\mathbf{z} = \begin{bmatrix} y \\ x \end{bmatrix}$$

$$\mu_z = \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i$$

$$C_z = \frac{1}{N} \sum_{i=1}^N (\mathbf{z}_i - \mu_z)(\mathbf{z}_i - \mu_z)^T$$

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}$$

$$C_z = \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix}$$

44

11755/18797

23 Oct 2012

## Learning the parameters of the Gaussian

$$\mu_z = \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i$$

$$C_z = \frac{1}{N} \sum_{i=1}^N (\mathbf{z}_i - \mu_z)(\mathbf{z}_i - \mu_z)^T$$

$$\mu_z = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}$$

$$C_z = \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix}$$

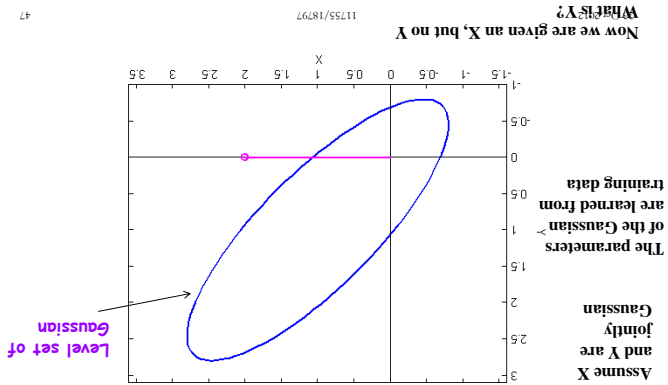
$$\mu_x = \frac{1}{N} \sum_{i=1}^N x_i$$

$$C_{XY} = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$$

45

11755/18797

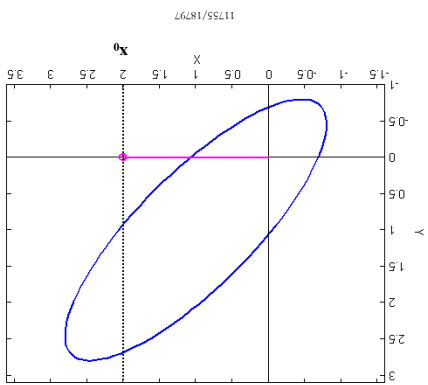
## MAP Estimator for Gaussian RV



47

11755/18797

## MAP estimator for Gaussian RV

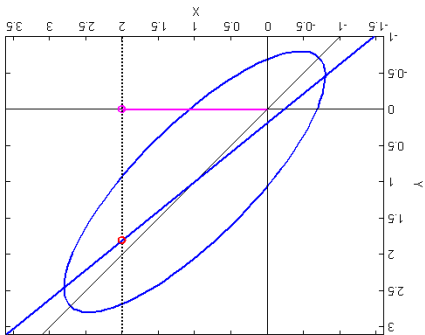


48

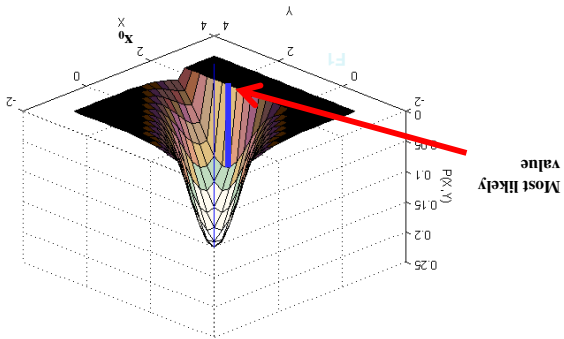
11755/18797

23 Oct 2012

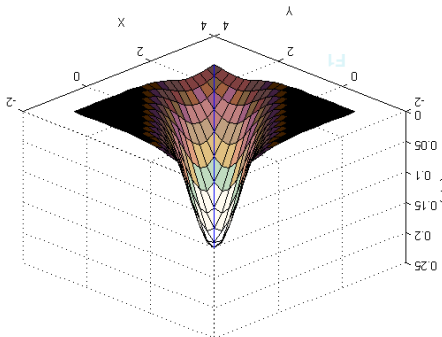




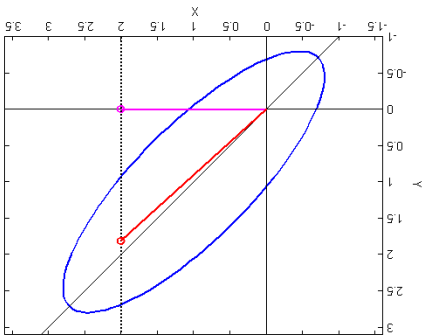
MAP Estimation of a Gaussian RV



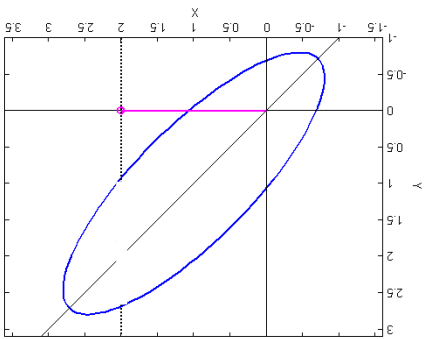
MAP estimation: The Gaussian at a particular value of X



MAP estimation: Gaussian PDF

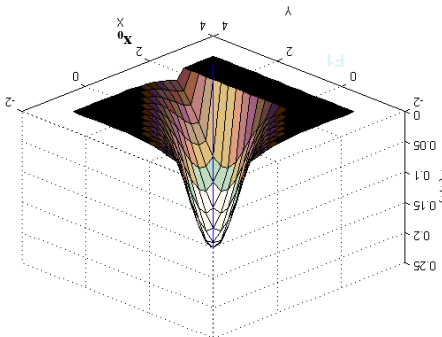


MAP Estimation of a Gaussian RV



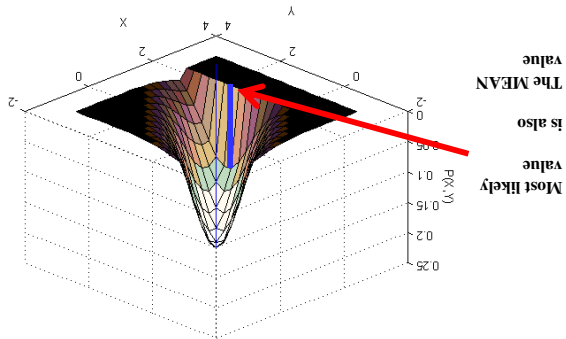
$Y = \text{argmax}_y P(y|X) ???$

MAP Estimation of a Gaussian RV



MAP estimation: The Gaussian at a particular value of X

- Would be true of any symmetric distribution



For the Gaussian: MAP = MMSE

- General principle of MMSE estimation:
  - y is unknown, x is known
  - Must estimate it such that the expected squared error is minimized
- Minimize above term

$$Err = E[\|y - \hat{y}\|^2 | x]$$

Its also a minimum-mean-squared error estimate

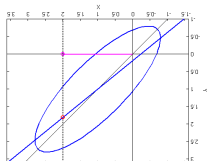
- Derivation? Later in the program a bit

$$\hat{y} = \mu_x + C_{yx}C_{xx}^{-1}(x - \mu_x)$$

- Scalar version given; vector version is identical

$$\hat{y} = \mu_x + C_{yx}C_{xx}^{-1}(x - \mu_x)$$

- Clearly a line
- Equation of line:



So what is this value?

- Just a weighted combination of the MMSE estimates from the component distributions

$$= \sum_k P(k)E[y | k, x]$$

$$E[y | x] = \int y \sum_k P(k)P(y | k, x)dy = \sum_k P(k) \int y P(y | k, x)dy$$

- Let P(y|x) be a mixture density
- The MMSE estimate of y is given by

$$P(y | x) = \sum_k P(k)P(y | k, x)$$

distributions

MMSE estimates for mixture

$$\hat{y} = E[y | x]$$

The MMSE estimate is the mean of the distribution

$$dErr = 2E[y^T y + y^T \hat{y} - 2y^T y | x] = 2\hat{y}^T d\hat{y} - 2E[y | x]^T d\hat{y} = 0$$

- Differentiating and equating to 0:

$$Err = E[y^T y + y^T \hat{y} - 2y^T y | x] = E[y^T y | x] + \hat{y}^T E[y | x] - 2\hat{y}^T E[y | x]$$

$$Err = E[\|y - \hat{y}\|^2 | x] = E[(y - \hat{y})^T (y - \hat{y}) | x]$$

- Minimize error:

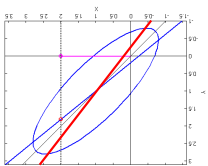
estimate

Its also a minimum-mean-squared error

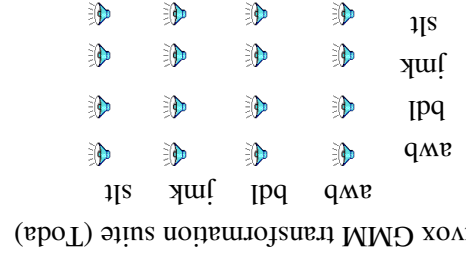
- What about the ML estimate of y
- Again, ML estimate of y, not regression parameter

- This is the MAP estimate of y
- NOT the regression parameter

$$\hat{y} = \mu_x + C_{yx}C_{xx}^{-1}(x - \mu_x)$$



This is a multiple regression



- Festvox GMM transformation suite (Toda)

## M MSE with GMM: Voice Transformation

$$E[y | x] = \sum_k P(k | x) (\mu_{k,y} + C_{k,yx} C_{k,xx}^{-1} (x - \mu_{k,x}))$$

$$E[y | x] = \sum_k P(k | x) E[y | k, x]$$

$$P(y | x) = \sum_k P(k | x) N(y; \mu_{k,y} + C_{k,yx} C_{k,xx}^{-1} (x - \mu_{k,x}), \Theta)$$

- $E(y|X)$  is also a mixture
- $P(y|x)$  is a mixture density

## M MSE estimates from a Gaussian mixture

$$P(y | x) = \sum_k P(k | x) P(y | x, k)$$

$$P(y | x) = \frac{P(x, y)}{P(x)} = \frac{P(x)}{P(x)} = \frac{\sum_k P(k, x, y)}{\sum_k P(k, x) P(y | x, k) P(x)}$$

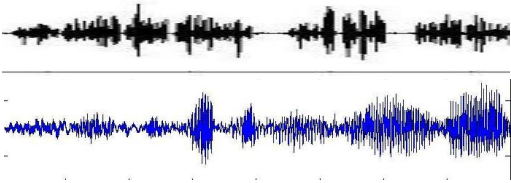
- Let  $P(y|x)$  is also a Gaussian mixture

$$P(x, y) = P(z) = \sum_k P(k) N(z; \mu_k, \Sigma_k) \quad z = \begin{bmatrix} y \\ x \end{bmatrix}$$

- Let  $P(x,y)$  be a Gaussian Mixture

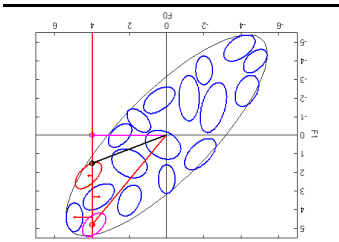
## M MSE estimates from a Gaussian mixture

- Align training recordings from both speakers
- Cepstral vector sequence
- Learn a GMM on joint vectors
- Given speech from one speaker, find MMSE estimate of the other
- Synthesize from cepstra



## Voice Morphing

- A mixture of estimates from individual Gaussians



## mixture

## M MSE estimates from a Gaussian mixture

$$P(y | x) = \sum_k P(k | x) N(y; \mu_{k,y} + C_{k,yx} C_{k,xx}^{-1} (x - \mu_{k,x}), \Theta)$$

$$P(y | x, k) = N(y; \mu_{k,y} + C_{k,yx} C_{k,xx}^{-1} (x - \mu_{k,x}), \Theta)$$

$$P(y, x, k) = N([y; x]; [\mu_{k,y}; \mu_{k,x}], \begin{bmatrix} C_{k,yy} & C_{k,yx} \\ C_{k,xy} & C_{k,xx} \end{bmatrix})$$

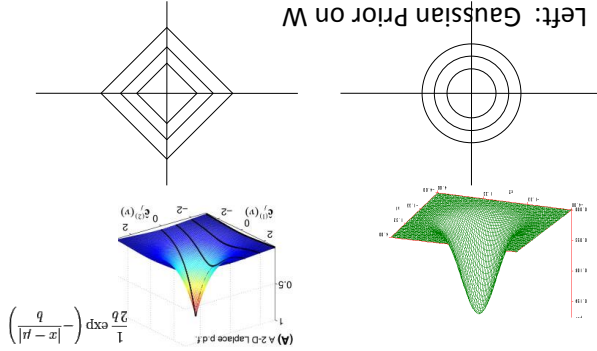
$$P(y | x) = \sum_k P(k | x) P(y | x, k)$$

- Let  $P(y|x)$  is a Gaussian Mixture

## mixture

## M MSE estimates from a Gaussian mixture

- Left: Gaussian Prior on W
- Right: Laplacian Prior



MAP estimate priors

- Similar to ML estimate with an additional term

$$\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} C^T - \log \sigma - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{a}^T \mathbf{X})^T (\mathbf{y} - \mathbf{a}^T \mathbf{X}) - 0.5\sigma^2 \mathbf{a}^T \mathbf{a}$$

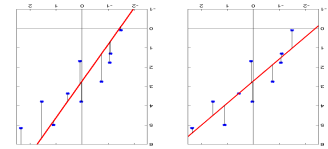
$$\log P(\mathbf{y} | \mathbf{X}, \mathbf{a}) = C - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{a}^T \mathbf{X})^T (\mathbf{y} - \mathbf{a}^T \mathbf{X})$$

- $P(\mathbf{a}) = N(0, \sigma^2 I)$
- $\log P(\mathbf{a}) = C - \log \sigma - 0.5\sigma^2 \|\mathbf{a}\|^2$

$$\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} \log P(\mathbf{a} | \mathbf{y}, \mathbf{X}) = \arg \max_{\mathbf{a}} \log P(\mathbf{y} | \mathbf{X}, \mathbf{a}) P(\mathbf{a})$$

MAP estimation of weights

- ML fit is sensitive
  - Small variations in data  $\rightarrow$  large variations in weights
  - Error is squared
  - Outliers affect it adversely
- Unstable
  - If dimension of  $X >=$  no. of instances
  - $(XX^T)$  is not invertible



$$\mathbf{A} = (\mathbf{X}\mathbf{X}^T)^{-1} \mathbf{X}\mathbf{y}^T$$

A problem with regressions

- Equivalent to *diagonal loading* of correlation matrix
  - Improves condition number of correlation matrix
  - Can be inverted with greater stability
  - Will not affect the estimation from well-conditioned data
  - Also called Tikhonov Regularization
  - Dual form: Ridge regression
- MAP estimate of weights
  - Not to be confused with MAP estimate of  $\mathbf{y}$

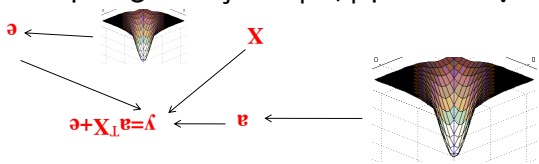
$$\mathbf{a} = (\mathbf{X}\mathbf{X}^T + \sigma I)^{-1} \mathbf{X}\mathbf{y}^T$$

$$dL = (2\mathbf{a}^T \mathbf{X}\mathbf{X}^T + 2\mathbf{y}\mathbf{X}^T + 2\sigma I) d\mathbf{a} = 0$$

MAP estimation of weights

$$\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} \log P(\mathbf{a} | \mathbf{y}, \mathbf{X}) = \arg \max_{\mathbf{a}} \log P(\mathbf{y} | \mathbf{X}, \mathbf{a}) P(\mathbf{a})$$

- Maximum *a posteriori* estimate
  - $\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} \log P(\mathbf{y} | \mathbf{X}, \mathbf{a})$
- Max. Likelihood estimate
  - $P(\mathbf{a}) = N(0, \sigma^2 I)$
- Assume weights drawn from a Gaussian



MAP estimation of weights

- No closed form solution
  - Quadratic programming solution required
  - Non-trivial

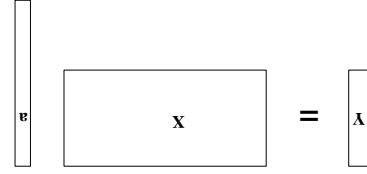
$$\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} C^T - (\mathbf{y} - \mathbf{a}^T \mathbf{X})^T (\mathbf{y} - \mathbf{a}^T \mathbf{X}) - \lambda^{-1} \|\mathbf{a}\|$$

- Assume weights drawn from a Laplacian
  - $P(\mathbf{a}) = \lambda^{-1} \exp(-\lambda^{-1} \|\mathbf{a}\|)$
  - Maximum *a posteriori* estimate

Laplacian prior

MAP estimation of weights with

## LASSO and Compressive Sensing



- Given  $Y$  and  $X$ , estimate sparse  $W$
- LASSO:
  - $X$  = explanatory variable
  - $Y$  = dependent variable
  - $a$  = weights of regression
- CS:
  - $X$  = measurement matrix
  - $Y$  = measurement
  - $a$  = data

- Various convex optimization algorithms
- LARS: Least angle regression
- Pathwise coordinate descent..
- Matlab code available from web

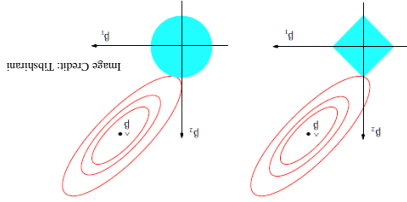
## LASSO Algorithms

- Identical to L1 regularized least-squares estimation
- $\hat{a} = \arg \max_a C^T (y - a^T X)^T (y - a^T X)^T - \lambda^{-1} \|a\|_1$
- Maximum *a posteriori* estimate
  - $P(a) = \lambda^{-1} \exp(-\lambda^{-1} \|a\|_1)$
- Assume weights drawn from a Laplacian

## MAP estimation of weights with Laplacian prior

## L1-regularized LSE

- Dual formulation
  - Quadratic programming solutions required
  - No closed form solution
- "LASSO" – Least absolute shrinkage and selection operator
  - $\hat{a} = \arg \max_a C^T (y - a^T X)^T (y - a^T X)^T - \lambda^{-1} \|a\|_1$



## Regularized least squares

- Regularization results in selection of suboptimal (in least-squares sense) solution
  - One of the loci outside center
- Tikhonov regularization selects *shortest* solution
- L1 regularization selects *sparsest* solution

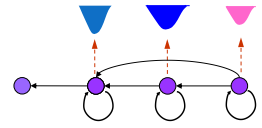
## An interesting problem: Predicting War!

- Economists measure a number of social indicators for countries weekly
  - Happiness index
  - Hunger index
  - Freedom index
  - Twitter records
  - ...
- Question: Will there be a revolution or war next week?

## An interesting problem: Predicting War!

- Issues:
  - Dissatisfaction builds up – not an instantaneous phenomenon
  - Usually
  - War / rebellion build up much faster
  - Often in hours
- Important to predict
  - Preparedness for security
  - Economic impact

## A Step Aside: Predicting Time Series



- An HMM is a model for time-series data
- How can we use it predict the future?

## Predicting with an HMM

$$P(s_t = s | O_{1:t}) = P(O_{1:t} = s, O_{1:t}) / \sum_{s'} P(O_{1:t} = s', O_{1:t})$$

- Given  $O_{1:t}$ 
  - Compute  $P(O_{1:t} = s)$
  - Using the forward algorithm – computes  $\alpha(s, t)$

$$\alpha(s, t) = P(O_1, O_2, \dots, O_t, \text{state}(t) = s)$$

## Predicting War



- Predict probability of unrest next week
  - Sequence of economic indicators for each week
  - Sequence of unrest markers for each week
  - At the end of each week we know if war happened or not that week
- This could be a new unrest or persistence of a current one

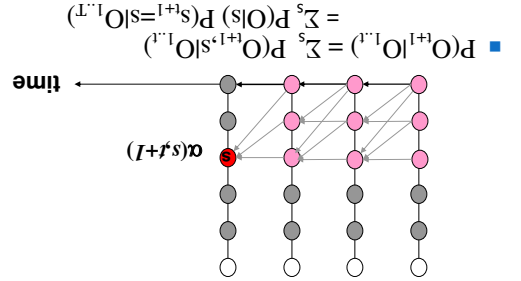
## Predicting with an HMM

- Given
  - Observations  $O_1 \dots O_t$
  - All HMM parameters
  - Learned from some training data
- Must estimate future observation  $O_{t+1}$ 
  - Estimate must consider *entire* history ( $O_{1:t}$ )
  - No knowledge of actual state of the process at any time

## Predicting with an HMM

- Given  $P(s^t = s | O_{1:t})$  for all  $s$ 
  - $P(s^{t+1} = s | O_{1:t}) = \sum_{s'} P(s^t = s' | O_{1:t}) P(s^t = s | O_{1:t}, P(s^t = s'))$
  - $P(O_{t+1} = s | O_{1:t}) = P(O_{1:t}) P(s^{t+1} = s | O_{1:t})$
  - $P(O_{t+1} = s) = \sum_{s'} P(O_{1:t}) P(s^{t+1} = s | O_{1:t})$
- This is a mixture distribution

## Predicting with an HMM

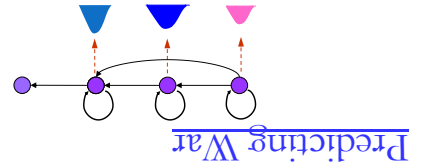


- $P(O_{t+1}|O_{1..t}) = \sum_s P(O_{t+1}, s|O_{1..t})$
- MIMSE estimate of  $O_{t+1}$  given  $O_{1..t}$ 
  - $E[O_{t+1} | O_{1..t}] = \sum_s P(s^{t+1}=s|O_{1..t}) E[O|s]$
- A weighted sum of the state means

23 Oct 2012

11755/18797

85



- Train an HMM on  $z = [w, s]$
- After the  $t^{\text{th}}$  week, predict probability distribution:
  - $P(z_t | z_1 \dots z_t) = P(w, z | z_1 \dots z_t)$
- Marginalize out  $x$  (not known for next week)
  - $P(w | z_1 \dots z_t) = \int P(w, s | z_1 \dots z_t) ds$
- War?  $\rightarrow E[w | z_1 \dots z_t]$

23 Oct 2012

11755/18797

87

## Predicting with an HMM

- MIMSE Estimate of  $O_{t+1} = E[O_{t+1} | O_{1..t}]$ 
  - $E[O_{t+1} | O_{1..t}] = \sum_s P(s^{t+1}=s|O_{1..t}) E[O|s]$
- If  $P(O|s)$  is a GMM
  - $E(O|s) = \sum_k P(k|s) \mu_{k,s}$

$$\hat{O}_{t+1} = \sum_k P(s|O_{1..t}) \sum_{k,s} w_{k,s} \mu_{k,s}$$

$$\hat{O}_{t+1} = \sum_s \frac{\alpha(t,s)}{\sum_k \alpha(t,s)} \sum_{k,s} w_{k,s} \mu_{k,s}$$

23 Oct 2012

11755/18797

86