

# SPARSE LINEAR MODELS

SELECTING SETS OF VARS WITH A MODEL-BASED APPROACH. IF LINEAR MODEL  $P(y|x) = P(y|f(W^T x))$  WE DO IT BY PROMOTING  $W$  TO BE SPARSE  
 → LOTS OF ZEROS USEFUL FOR: SMALL  $N$ , LARGE  $D$ ; KERNEL MACHINES; SPARSE WAVELET REPRESENTATION FOR SIGNALS.

## BAYESIAN VARIABLE SELECTION

$$P(y|D) = \frac{e^{-f(y)}}{\sum_{y'} e^{-f(y')}} \quad \text{POSTERIOR OVER MODELS} \quad f(y) = -[\log P(D|y) + \log P(y)] \quad \text{MAP: } \hat{y} = \text{ARGMAX}_y P(y|D) = \text{ARGMIN}_y f(y)$$

$$\text{MEDIAN MODEL: } \hat{y} = \{y: P(y, = 1|D) > 0.5\}$$

MARGINAL INCLUSION PROBS

$$\text{SPRUE AND SLAB MODEL} \quad \log P(y|D) \approx \log P(y|x, \hat{w}, \hat{\sigma}^2) - \frac{\|x\|_0}{2} \log N - \lambda \|x\|_0 + \text{CONST}$$

•  $\|x\|_0 = \ell_0$  NORM, NUM OF NON-NULL ELEMENTS

BIC APPROX  
PENALTY

PRIOR  
PENALTY

## BEER/GAUSS MODEL

$$\ell_0\text{-REGULARIZATION: } f(w) = \|y - Xw\|_2^2 + \lambda \|w\|_0 \rightarrow \text{CONVERTS DISCRETE OPTIMIZATION PROBLEM OVER } \gamma \text{ INTO A CONTINUOUS ONE OVER } w \in \mathbb{R}^D$$

MAP      STILL NONSMOOTH AND HARD 2 OPTIMIZE

## OPTIMIZATION ALGOS

SPACE OF MODELS IS  $2^D \rightarrow$  WE NEED HEURISTICS. FIT THE MODEL AT EACH POINT. COMPUTING  $P(D|\gamma)$  OR  $\int P(D|w)P(w)dw \rightarrow$  MAXIM  $\rightarrow$  WRAPPED METHOD  
 FOR EFFICIENCY WE WANT TO BE ABLE TO COMPUTE  $\gamma'$  GIVEN  $\gamma \rightarrow$  UPSAMPLE SUFF. STATS.  $\rightarrow$  OR IF  $\gamma'$  DIFFERS BY 1 BIT AM  $f(\gamma)$  DEPENDS ON DATA ONLY VIA  $X$

## GREEDY SEARCH

- SINGLE BEST REPLACEMENT  $\rightarrow$  GREEDY HILL CLIMBING BY MOVING IN NEIGHBORHOOD OF  $\gamma$  (1 BIT FLIPS). LATTICE SEARCH.
- ORTHOGONAL LEAST SQUARES  $\rightarrow$  START WITH  $\lambda=0$  AND EMPTY SET. ADD BEST FEATURE AT EACH STEP.  $\gamma^* = \text{ARG MIN}_\gamma \min_w \|y - (X_\gamma w)\|_2^2$  D-0th LS PER STEP
- ORTHOGONAL MATCHING PURSUITS  $\rightarrow$  OLS IS EXPENSIVE. FREEZE WEIGHT AT CUR VAL THEN PICK NEXT  $\gamma^* = \text{ARG MIN}_\gamma \min_w \|y - Xw - Bx_\gamma\|_2^2$  ONLY 1 LS PER STEP BUT NOT AS ACCURATE
- MATCHING PURSUITS:  $\rightarrow$  VERY GREEDY, ONLY ADDS COLUMN MOST CORRELATED W/ CURRENT RESIDUAL | LS BOOSTING
- BACKWARDS SELECTION:  $\rightarrow$  STARTS WITH ALL VARS, REMOVES WORST ONE AT EACH STEP
- F-B1: SBR WITH OMP STEP FOR NEXT SELECTION
- BAYESIAN MATCHING PURSUITS: OMP BUT USES MARGINAL LIKELIHOOD CRITERION VS LEAST SQUARES. ALSO BEAM SEARCH.

## STOCHASTIC SEARCH

APPROXIMATES POSTERIOR, MCMC. BUT INEFFICIENT. USE  $P(y|D) = \frac{e^{-f(y)}}{\sum_{y'} e^{-f(y')}}$

## EM/VARIATIONAL

EM ON SPRUE/SIAB, WITH NARROW GAUSSIAN INSTEAD OF DELTA. LOCAL MINIMA ISSUE.

EM ON BEER/GAUSS  $\rightarrow$  ALL VAR TRAINED AWAY. USE MEAN FIELD APPROXIMATION

# L1 REGULARIZATION

PUTTING A 0 MEAN LAPLACE PRIOR ON THE PARAMS AND PERFORMING MAP. CAN BE COMBINED WITH ANY CONVEX/NO CONVEX NLL.

IT ENCOURAGES  $w_j = 0$

PRIOR:  $P(w|\lambda) = \prod \text{LAP}(w_j | 0, 1/\lambda) \propto \prod e^{-\lambda |w_j|}$

NLL:  $f(w) = -\log P(D|w) - \log P(w|\lambda) = \text{NLL} + \lambda \|w\|_1$

LIN REG  
LOGISTIC REG  
GLM  
...

$\|w\|_1 = \sum |w_j|$  L1 NORM

CONVEX APPROX OF L0 OBJECTIVE ARGMIN  $\text{NLL}(w) + \lambda \|w\|_1$

IN LINEAR REGRESSION:  $f(w) = \sum -\frac{1}{2\sigma^2} (y_i - (w_0 + w^T x))^2 + \lambda \|w\|_1 = \text{RSS}(w) + \lambda' \|w\|_1, \lambda' = 2\sigma^2 \lambda$

YIELDS SPARSE SOLUTIONS

$\ell_1$  MIN  $\text{RSS}(w) + \lambda \|w\|_1 \rightarrow \text{MIN } \text{RSS}(w) \text{ ST } \|w\|_1 \leq B$  LASSO EQUATION

NONSMOOTH

SMOOTH BUT CONSTRAINED



- MORE PROBABLE THAT CORNERS INTERSECT ELLIPSE  $\rightarrow$  CORNERS ARE SPARSE SOLUTIONS
- SPARSE SOLUTIONS COST LESS

$\ell_2$  MIN  $\text{RSS}(w) \text{ ST } \|w\|_2^2 \leq B$



- NO CORNERS, NO PREFERENCE FOR SPARSITY
- COST,  $\ell_2$  NORM, FOR SPARSE IS SAME AS FOR DENSE

$$\begin{cases} (c_j + \lambda)/a_j & c_j \leq -\lambda \\ 0 & -\lambda \leq c_j \leq \lambda \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases}$$

- SETS TO 0 OR SHRINKS
- SOFT THRESHOLDING

$\lambda_{\text{max}} = \max_j |v_j \text{NLL}(0)|$ ; IF  $\lambda > \lambda_{\text{max}} \rightarrow \hat{w} = 0$

## COMPARISON

- MLE  $\hat{w}_{\text{OLS}} = X^T y$
- LASSO  $w_{\text{LASSO}} = \text{SIGN}(\hat{w}_{\text{OLS}}) (|\hat{w}_{\text{OLS}}| - \frac{\lambda}{2})$
- RIDGE  $\hat{w}_{\text{RIDGE}} = \frac{\hat{w}_{\text{OLS}}}{1 + \lambda}$
- SUBSET SELECTION  $\begin{cases} w_{\text{OLS}} & \text{IF RANK} \leq K \\ 0 & \text{ELSE} \end{cases}$

## REGULARIZATION PATH

- PLOTS  $\hat{w}_j(\lambda)$  VS  $\lambda$  FOR EACH FEATURE  $j$ . SPARSE SOLUTIONS FOR BOUND  $\ell_1$  BETWEEN 0 AND  $B_{\text{max}} = \|w_{\text{OLS}}\|_1$ . IN BETWEEN EACH LINEAR COEFF INCREASES OR DECREASES LINEARLY. PIECEWISE LINEAR. LARS ALGORITHM; SAME COST AS LEAST SQUARES FIT  $O(\min\{ND^2, DN^2\})$
- LIMITATION IF  $D > N$  LASSO ONLY GETS TO  $N$  BEFORE GOING TO COMPLETE SET FOR OPTIMAL OLS  $\rightarrow$  USE ELASTIC NET

## MODEL SELECTION

- RECOVER TRUE SPARSITY PATTERN: PICK  $\lambda$ . USE CV.  $\ell_1$  WHEN IF PROVIDED DATA. DO BOOTSTRAP AND COMPARE HOW EACH VAR IS PICKED IN RUNS AND OBTAIN SPARSE ESTIMATOR  $\rightarrow$  BODASSO PICK IF VAR RETURNED AT LEAST 90% OF TIMES FOR GIVEN  $\lambda$
- OPTIMAL  $\lambda$  FOR PREDICTION  $\rightarrow$  MODEL SELECTION CONSISTENCY

BAYESIAN INFERENCE: ONLY POSTERIOR MODE IS SPARSE. POSTERIOR MEAN BETTER IF WE MINIMIZE SQUARED ERROR.  $P(\text{MEAN} + \text{SPRUE AM SUB}) > P(\text{MODE} + \text{LAP PRIOR})$

## L1 REGULARIZATION ALGORITHMS

- COORDINATE DESCENT: OPTIMIZE VARS 1-04-1  $w_j = \text{ARGMIN}_j (f(w + ze_j) - f(w))$ . PICK RANDOM OR ONE WHERE STEEPEST GRAD. OR IF ANALYTICAL SOLUTION FOR EACH 1D PROBLEM.  $\rightarrow$  SHOOTING ALGO
- LARS: ACTIVE SET  $\rightarrow$  MORE THAN 1 VAR AT TIME. HOMOTOPY  $\rightarrow$  STARTS FROM  $\lambda_{\text{max}}$  DOWN TO DESIRED  $\lambda_x$ . LEAST ANGLE REGRESSION & SHANNING. STAYS AT  $\lambda$  WHERE THE SINGLE MOST CORRELATED VAR W/Y IS PICKED.  $\lambda$  DECREASES UNTIL SECOND VAR WITH SAME CORRELATION WITH RESIDUAL FROM FIRST IS FOUND. NEW  $\lambda$  FOUND ANALYTICALLY, GEOMETRICALLY. DISCREPANT. ANALYTICAL SOLUTIONS ONLY FOR LINEAR, NOT OTHERS GLM.
- PROXIMAL & GRADIENT PROJECTION: GOOD FOR VERY LARGE PROBLEMS, CAN USE FOR OTHER TYPES BEYOND  $\ell_1$  REG. CONVEX OBJECTIVE  $f(w) = L(w) + R(w) \rightarrow f(w) = R(w) + \frac{1}{2} \|w - y\|_2^2$ . WE CALL THE MINIMIZER OF THIS PROXIMAL OPERATOR (DEPENDS ON  $R(w)$ ) IN MOST CASES  $\text{PROXR}(w)$  IS PROJECTION ON SET  $C = \{w: \|w\|_1 \leq B\} = \text{PROXR}(0) \rightarrow$  COMPUTED IN  $O(D)$  TIME. USED IN GRADIENT DESCENT ROUTINES QUADRATIC APPROX OF LOSS CENTERED ON  $\theta_n$ 
  - $\theta_{n+1} = \text{PROXR}(\theta_n - t_n g_n)$  APPROXIMATES LN POINT OTHER THAN  $\theta_n$
  - $\theta_n = \theta_n - t_n g_n$  NESTEROV METHOD:  $\begin{cases} \theta_{n+1} = \text{PROXR}(\theta_n - t_n g_n) \\ g_n = \nabla L(\theta_n) \\ \phi_n = \theta_n + \frac{k+1}{k+2} (\theta_n - \theta_{n-1}) \end{cases}$  NESTEROV + ITERATIVE SOFT THRESHOLDING + CONTINUATION = FISTA
  - $g_n = \nabla L(\theta_n)$  FAST ITERATIVE SHANNING THRESHOLDING

# PROXIMAL & PROJECTED GRADIENT METHODS

$f(\theta) = L(\theta) + R(\theta)$  -  $L$  LOSS, CONVEX AND DIFFERENTIABLE -  $R$  REGULARIZER NOT NECESSARILY DIFFERENTIABLE

• PROXIMAL OPERATOR FOR CONVEX FCN  $R$   $\text{PROX}_R(y) = \underset{z}{\text{ARGMIN}} \left( R(z) + \frac{1}{2} \|z - y\|_2^2 \right)$

- 'MOVES' FCN POINT TOWARDS MINIMUM OF FCN  $\text{PROX}_R(y) \approx y - \lambda \nabla R(y)$
- INTUITIVELY CONNECTED TO GRADIENT.
- FIXED POINT FOR PROX ARE FCN MINIMA

• EXAMPLE:  $L(\theta) = \text{RSS}$ ,  $R(\theta) = \|\theta\|_1$  = LASSO

-  $R(\theta) = \lambda \|\theta\|_1 \rightarrow \text{PROX} = \text{SOFT THRESHOLDING}$

-  $R(\theta) = \lambda \|\theta\|_0 \rightarrow \text{PROX} = \text{HARD THRESHOLDING}$

-  $R(\theta) = \lambda \|\theta\|_2 \rightarrow \text{PROX} = \text{PROX}_L(\theta) = \underset{z}{\text{ARGMIN}} \|z - \theta\|_2^2$

PROX IS EASIER TO COMPUTE

## PROJECTED PROXIMAL GRADIENT METHOD

-  $\theta_{n+1} = \underset{u}{\text{ARGMIN}} \left[ t_n R(z) + \frac{1}{2} \|z - u\|_2^2 \right] = \text{PROX}_{t_n R}(u_n)$ ,  $u_n = \theta_n - t_n g_n$ ,  $g_n = \nabla L(\theta_n)$ ,  $t_n = 1/\eta_n$ ,  $\eta_n$  = HESIAN APPROXIMATION

-  $R(\theta) = 0 \rightarrow$  VANILLA GRAD DESCENT

• CAN HAZ NESTEROV ACCELERATION

-  $R(\theta) = \|\theta\|_1 \rightarrow$  PROJECTED GRADIENT DESCENT

-  $R(\theta) = \lambda \|\theta\|_1 \rightarrow$  ITERATIVE SOFT THRESHOLDING  $\rightarrow$  ISTA  $\rightarrow$  USED FOR LASSO

## CONJUGATE GRADIENTS

$\rightarrow$  MAXIMAL STEP IN EACH DIMENSION, OPTIMAL  $\rightarrow$  ORTHOGONAL. IN PRACTICE: CONJUGATE/A-ORTHOGONAL CONJUGACY:  $p_i, p_j$  IFF  $p_i^T A p_j = 0$  FOR SOME MATRIX  $A$

ZEROS OF:  $\alpha_i = p_i^T \frac{(-\nabla f(x_{i-1}))}{p_i^T A p_i}$

HOW TO FIND CONJUGATE DIRECTIONS?

GRAMM-SCHMIDT PROCESS: TAKE CANDIDATE DIRECTION AND REMOVE PARTS IN DIR ALREADY DONE

$p_n = u_n + \sum_{i=0}^{n-1} \beta_{ni} p_i$  TWO OPTIONS FOR  $\beta_i$ : -  $\beta_{i+1} = \frac{\nabla f(x_{i+1})^T \nabla f(x_{i+1})}{\nabla f(x_i)^T \nabla f(x_i)}$  OR  $(\nabla f(x_{i+1}) - \beta \nabla f(x_i))^T \nabla f(x_i)$

• INITIAL DIRECTION  $p_0$  IS STEEPEST DESCENT

• FIND  $\alpha$  MINIMIZING  $f(x_i + \alpha p_i) \rightarrow x_{i+1} = x_i + \alpha p_i$

NEWTON - RHAPSON

$$\alpha = \frac{\nabla f(x)^T p}{p^T H(x) p}$$

•  $p_{i+1} = -\nabla f(x_{i+1}) + \beta_{i+1} p_i$

• RUN MULTIPLE TIMES ARE MUCH BETTER BECAUSE FINDS DIFFERENT DIRECTIONS

• GOOD BECAUSE REQUIRES NO MATRIX INVERSIONS

FLEISCHER - REEVES / POWELL - CONJUGATE



## EM FOR LASSO

MODEL WITH MIXING DISTRIBUTION ON VARIANCES

WE CAN REPRESENT VARIANCE DISTRIBUTION AS GAUSSIAN SCALE MIXTURE, DERIVE REPRESENTATION OF LASSO.

$$P(y, w, \tau, \sigma^2 | x) = \exp(-w/\sigma^2) \prod \exp(-\frac{\tau^2}{2}\tau)$$

E STEP INFERS  $\tau^2$  AND  $\sigma^2$ , M STEP ESTIMATES  $w$

SOMETIMES DOES NOT WORK FOR NUMERICAL REASONS

### WHY EM?

- PROVIDES WAY TO DERIVE  $\ell_1$  REG. FOR OTHER MODELS
- ALLOWS FOR TRYING OUT OTHER PRIORS ON VARIANCES
- CAN COMPUTE FULL POSTERIOR AND NOT JUST MAP  $\rightarrow$  BAYESIAN LASSO

## L1 EXTENSIONS

### - GROUP LASSO

WHEN VECTOR OF WEIGHTS AND NOT SINGLE, FOR EACH VAR. IE MULTINOMIAL LOGISTIC, LINEAR REG. WITH CATEGORICAL TASKS, MULTITASK LEARNING

POSITIONS PARAM VECTOR INTO GROUPS

$J(w) = \text{NLL}(w) + \sum \lambda_g \|w_g\|_2 \rightarrow$  IN HERE TWO NORM RESULTS IN SPARSITY. ELSE USE  $\infty$  NORM  $\|w_g\|_\infty = \max(w_g)$ . NLL IS LEAST SQUARES  $\rightarrow$  LASSO.

CAN BE SEEN AS GSM  $\rightarrow$  VARIANCE GROUP TERM COMES FROM A GAMMA PRIOR

ALGOS: PROXIMAL GRADIENT DESCENT DECOMPOSED IN G, EM

### - FUSED LASSO

IF WE WANT NEIGHBORING COEFFICIENTS TO BE SIMILAR TO EACH OTHER AND SPARSE. LOCATION-BASED PENALTY. GRAPHS, IMAGES, VIDEOS, ETC.

SIMILAR TO CHAIN-STRUCTURED GAUSSIAN MARKOV FIELDS, RANDOM WALKS. SOLVES WITH EM.

$$J(w) = J(w; \lambda_1, \lambda_2) = \sum (y_i - w_i)^2 + \lambda_1 \sum |w_i| + \lambda_2 \sum |w_{i+1} - w_i|$$

### - ELASTIC NET

COMBINES LASSO AND RIDGE

IF GROUP OF VARS CORRELATED  $\rightarrow$  LASSO LIKELY PICKS ONE OF THEM

$D \gg N \rightarrow$  LASSO AT MOST  $N$  VARS

$N > D$  + CORRELATED VARS  $\rightarrow$  RIDGE DOES BETTER

ALGORITHM! LASSO ON MODIFIED DATA  $\tilde{X} = \begin{pmatrix} X \\ \sqrt{\lambda_2/\lambda_1} I_D \end{pmatrix}$   $\tilde{y} = \begin{pmatrix} y \\ 0 \end{pmatrix}$   $C = (1 + \lambda_2)^{-1/2}$

USE LARS  $\rightarrow$  LARS-EN

IF WE STOP AT  $m$  VARS  $\rightarrow O(m^2 + Dm^2)$

CV TO PICK  $\lambda_1, \lambda_2$

$$\text{OBJECTIVE: } J(w; \lambda_1, \lambda_2) = \|y - Xw\|^2 + \lambda_2 \|w\|_2^2 + \lambda_1 \|w\|_1$$

STRICTLY CONVEX

GROUPING EFFECT =

HIGHLY CORRELATED VARS HAVE EQUAL WEIGHTS

$$\tilde{w} = \text{ARGMIN} \|\tilde{y} - \tilde{X}\tilde{w}\|_2^2 + C\lambda_1 \|\tilde{w}\|_1$$

$$w = C\tilde{w}$$

OBS! DOUBLE SHRINKAGE: RESCALE UP, UNDO  $\ell_2$ ,  $\hat{w} = \sqrt{1 + \lambda_2} \tilde{w}$

GSM: WITH PRODUCT OF GAUSSIAN AND LAPLACE FAVORS MAP, MCML, VARIATIONAL BAYES.

## NONCONVEX REGULARIZERS

LAPLACE PRIOR NOT SUPPRESSES NOISE. TOO MUCH SHRINKAGE ON LARGE VALUE COEFFS  $\rightarrow$  INTRODUCE MORE FLEXIBLE PRIORS NO GLOBAL OPTIMUM ANYMORE BUT OUTPERFORM IN PREDICTIVE ACCURACY AND VARIABLE SELECTION

### BAYES REGRESSION

GENERALIZES  $\ell_1$

$$\hat{w} = \text{NLL}(w) + \lambda \sum_j |w_j|^b$$

MAP WITH EXPONENTIAL POWER DISTRIBUTION



$b=0 \rightarrow \ell_0 \rightarrow \text{BSS}$

CAN GENERATE BIMODALS

$b=1 \rightarrow \text{LAPLACE} \rightarrow \text{LASSO}$

$b=2 \rightarrow \text{GAUSSIAN} \rightarrow \text{RIDGE}$

### HIERARCHICAL ADAPTIVE LASSO

DIFFERENT PENALTY PARAMETER FOR EACH PARAMETER. PENALTIES MODELED AS RV COMING FROM CONJUGATE PRIORS  $\gamma_j \sim \text{IG}(\alpha, b)$ ,  $\tau_j^2 | \gamma_j \sim \text{Ga}(\lambda, \gamma_j^2/2)$

• OFTEN WORKS MUCH BETTER THAN  $\ell_1$ , FIT WITH EM

• STAFFISH  $\rightarrow$  MORE AGGRESSIVE SPARSITY

• APPROXIMATE; HARD-THRESHOLDING

E-STEP = SAME FOR LASSO

M-STEP = WEIGHTED LASSO PROBLEM  $\rightarrow$  USE LARS

• CAN USE WITH ALL SORTS OF COMBINATIONS OF  $P(\tau_j^2)$ ,  $P(\gamma_j)$ ,  $P(w_j)$  FOR DIFFERENT RESULTS

• SCALED MIXTURE OF LAPLACIANS

$$w_j | \tau_j^2 \sim N(0, \tau_j^2)$$

## AUTOMATIC RELEVANCE DETERMINATION / SPARSE BAYESIAN LEARNING

ARD

SAL

APPROACH BASED ON TYPE 2 ML, EMPIRICAL BAYES. • WE INTEGRATE OUT  $w$  AND MAXIMIZE MARGINAL LIKELIHOOD WRT  $\tau$ . USING EM OR  $\ell_1$  SCHEME

• IS POINT ESTIMATION

• VARIANCES ESTIMATED  $\rightarrow$  PLUG-IN TO COMPUTE WEIGHT POSTERIOR MEANS  $E[w | \hat{\tau}, D] \rightarrow$  SPARSENESS

• NON-FACTOREAL PRIOR

• 'PUNISHES' SOLUTIONS THAT WASTE PROBABILITY MASS

• WE INTEGRATE  $w$  OUT AND OPTIMIZE  $\alpha$ , MLE STANDARD DOES OPPOSITE. PARAMS  $w_j$  BECOME CORRELATED DUE TO EXPANDING AWAY.  $\rightarrow$  FOR  $\alpha$ , WE ESTIMATE INFO FROM ALL FEATURES

• PRIOR IS NON-FACTOREAL:  $P(w | \alpha)$  DEPENDS ON DATA  $D$  AND  $\sigma^2$

• NON-FACTOREAL OBJECTIVE HAS ALWAYS FEWER LOCAL MINIMA THAN FACTOREAL ONES  $\rightarrow$  BUT STILL GLOBAL MAXIMUM IS EQUAL TO  $\ell_0$

• ALGORITHMS: EM, FIXED-POINT ALGORITHM, REWEIGHTED  $\ell_1$   $w^{(t+1)} = \text{ARGMIN } \text{NLL}(w) + \sum_j \lambda_j |w_j|$ , VARIATIONAL APPROXIMATIONS

$\begin{matrix} \text{EM} \\ \text{MSE} \\ \alpha, \rho \end{matrix}$

• IF GAUSSIAN PROCESSES, NNLS, ...

$$K(x, x') = \theta_0 \exp \left[ -\frac{1}{2} \sum_i \eta_i (x_i - x'_i)^2 \right] \quad \eta \rightarrow 0: \text{FOR INSENSITIVE TO } x \text{ VALUES}$$

USE ML TO ESTIMATE  $\eta \rightarrow$  DESERT VARS WITH LITTLE POSTERIOR EFFECT, CAN DISCARD

## SPARSE CODING

IDEA: SPARSE PRIORS FOR UNSUPERVISED LEARNING, THINK ICA BUT WITH SPARSITY PROMOTING PRIORS. VECTORS  $x_i$  AS SPARSE COMBINATION OF BASIS VECTORS

SPARSE CODING:  $w$  NOT ORTHOGONAL

SPARSE PRIOR ON LATENT FACTORS

$w$ : HERE, CALLED A DICTIONARY COLUMN OF  $w$ : ATOM

• IF  $LTD$  OVERCOMPUTES

FIXED: WAVELET/DCT BASIS

LEARNED:

SPARSE PCA: SPARSE ON  $w$ s, PRIOR IS GAUSS

SPARSE MATRIX FACTORIZATION: SPARSE  $w$ s + SPARSE LF

$$\log P(D | w) = \sum_i \log \int N(x_i | w z_i, \sigma^2 I) P(z_i) dz_i$$

## HOW TO LEARN DICTIONARY!

$P(z)$

LAPLACE:  $\text{NLL}(w, z) = \sum_i \frac{1}{2} \|x_i - w z_i\|_2^2 + \lambda \|z_i\|_1$  CONSTRAIN  $\ell_2$  OF COLUMNS TO BE  $\leq 1$

• FOR FIXED  $z$ , LEAST SQUARE OPTIMIZATION OVER  $w$

• FOR FIXED  $w$  LASSO PROBLEM OVER  $z$

• OTHER MODELS  $\rightarrow$  OTHER OPTIMIZATION PROBLEMS

• SMF: ELASTIC NET TYPE PENALTY ON WEIGHTS =  $\min_{w, z} \frac{1}{2} \sum_i \|x_i - w z_i\|_2^2 + \lambda \|z\|_1$

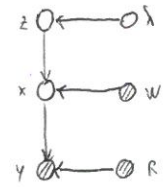
• ANALYSIS/SYNTHESIS LOOP: ALTERNATE OPTIMIZE  $w$  AND  $z$

• CAN USE OTHER SPARSITY INDUCING PRIORS THAN LAPLACE, I.E. BIMODALS, BETA PROCESSES

# WHAT TO DO WITH SPARSE CODING

• COOL BECAUSE LEARNED BASIS VECTOR LOOK LIKE BRAIN 'VISUAL FILTERS', ALSO ICA, BUT NOT PCA

• COMPRESSED SENSING:  $y = Rx + e$ , LOW-DIM PROJECTION OF DATA.  $R$  KNOWN SENSING MATRIX. IE MRI EACH BEAM DIRECTION IS ROW IN  $R$ .  
 → INFER  $P(x|y, R)$  VIA BAYESIAN INFERENCE. ASSUMING  $x = W \cdot z$   
 ↳ SPARSE PRIOR  
 ↳ DICTIONARY



## • IMAGE INPAINTING / DENOISING

WE PARTITION IMG IN OVERLAPPING PATCHES  $y_i$  AND CONCATENATE TO FORM  $y$ .  
 $R$  DEFINED TO SELECT PATCH  $i$  WITH  $i$ TH ROW.  $V$  VISIBLE,  $H$  HIDDEN.

WE OBSERVE LOW-DIM  $y$   
 PASSED THROUGH  $R$ .

WE COMPUTE  $P(y_H | y_V; \theta)$   $\theta$  ARE PARAMS  $W$  AND SPARSITY LEVEL  $\lambda$

DICTIONARY LEARNED OR FIXED.

$x$  IS SPARSE WRT  $W$  AND  $z$ .

ALTERNATIVE: FIELD OF EXPERTS → ENCODE CORRELATIONS BETWEEN NEIGHBORING PATCHES / NO LATENTS → COMPUTATIONALLY HEAVIER