Introduction to Machine Learning

Brown University CSCI 1950-F, Spring 2012 Prof. Erik Sudderth

Lecture 14: Sparsity & L₁ Regularization: The Lasso

> Many figures courtesy Kevin Murphy's textbook, Machine Learning: A Probabilistic Perspective

Feature Selection for Regression $f(w) = ||y - \Phi w||_{2}^{2} + \lambda ||w||_{0}$ $||w||_{p} = \left(\sum_{j=1}^{D} |w_{j}|^{p}\right)^{1/p} \qquad ||w||_{0} = \underset{nonzero \ entries}{nonzero \ entries}$

• L₀ penalty comes from Bernoulli prior on feature usage indicators, and large variance prior on non-zero weights:

$$p(\boldsymbol{\gamma}) = \prod_{j=1}^{D} \operatorname{Ber}(\gamma_j | \pi_0) = \pi_0^{||\boldsymbol{\gamma}||_0} (1 - \pi_0)^{D - ||\boldsymbol{\gamma}||_0} \qquad ||\boldsymbol{\gamma}||_0 = \sum_{j=1}^{D} \gamma_j$$
$$w_j \sim \mathcal{N}(0, \sigma_w^2) \qquad \sigma_w^2 \to \infty$$

- L_2 penalty comes from Gaussian regression likelihood: $p(y_i \mid x_i, w, \gamma, \sigma^2) = \mathcal{N}(y_i \mid \sum_{j=1}^{D} \gamma_j w_j \phi_j(x_i), \sigma^2)$
- **Problem:** Optimization is hard combinatorial problem!

Greedy Deterministic Search Backward Selection $\{1, 2, 3, 4\}$ $\{1, 2, 3\}$ $\{2, 3, 4\}$ $\{1, 3, 4\}$ $\{1, 2, 4\}$ $\{1,2\}$ $\{1,3\}$ $\{1,4\}$ $\{2,3\}$ $\{2,4\}$ $\{3,4\}$ $\{1\}$ $\{2\}$ $\{3\}$ $\{4\}$ {}

Forward Selection

- Consider all possible ways of adding *(forward selection)* or removing *(backward selection)* one feature
- Add or remove the best feature, or stop if the current model is best
- Wrapper method: Can be applied to any objective. *Guarantees???*



Pascal's Triangle (http://www.mathwarehouse.com/)



When used as a prior on vectors of model parameters:

- Compared to Gaussian, stronger bias that many near zero
- When find MAP estimate, some weights are *exactly* zero
- Learning harder than for Gaussian, but still *convex*

Constrained Optimization



Where do level sets of the quadratic regression cost function first intersect the constraint set?

Gradient-Based Optimization

Laplacian prior L₁ regularization Lasso regression Gaussian prior L₂ regularization Ridge regression



Generalized Norms: Bridge Regression



- Convex objective function (true norm): $b \ge 1$
- Encourages sparse solutions (cusp at zero): $b \le 1$
- Lasso/Laplacian (convex & sparsifying): b = 1
- Ridge/Gaussian (classical, closed form solutions): *b* = 2
- Sparsity via discrete counts (greedy search): $b \rightarrow 0$

Bayesian Linear Regression

0 data points observed



Bayesian Linear Regression

1 data point observed



Bayesian Linear Regression

2 data points observed





Shrinkage for Orthonormal Features



feature selection

remain non-zero

Regularization Paths

Prostate Cancer Dataset with N=67, D=8



Vertical lines are models chosen by cross-validation

Optimization: Projected Gradient



Generic method based on gradient & projection operators: $\mathbf{w}_{k+1} = \mathbf{w}_k + \eta_k \mathbf{d}_k$ $\mathbf{d}_k = \operatorname{proj}(\mathbf{w}_k - \eta_k \mathbf{g}_k) - \mathbf{w}_k$

Projection onto non-negativity constraint is trivial: $w_i := \max(w_i, 0)$

Good properties, extensions choose even better descent directions...



- MAP estimate is atypical: Samples from the prior have all weights non-zero with probability one
 - Full Bayesian learning requires richer prior models (often perform better at greater computational cost)
- In addition to setting some coefficients exactly to zero, non-zero coefficients are significantly biased towards zero
 - Two-stage estimators estimate weight vector support, then re-estimate weights with a less strong prior
- Theory guarantees *support recovery* in some conditions, but can be *unstable* when features are strongly correlated