

Probabilistic Models for Sparse Regression and Classification

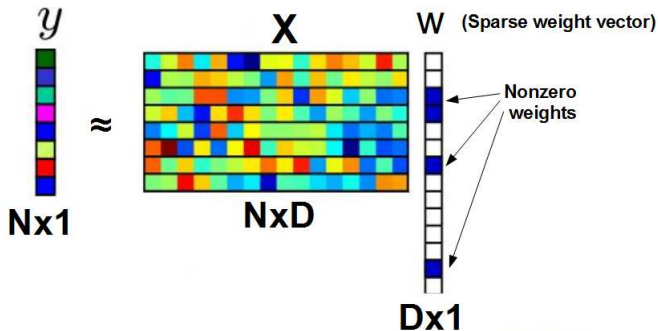
Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

March 13, 2018

Sparse Linear Models

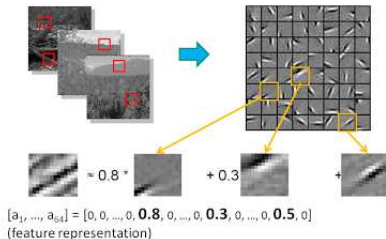
- Consider a linear (or generalized linear) model: $\mathbf{y} \approx \mathbf{X}\mathbf{w}$
- D features in the data. Suppose very few are relevant for predicting \mathbf{y}
- We would expect \mathbf{w} to be very sparse (very few nonzeros in it)



- Very general model, appears not just in regression/classification, but also in many other problems (e.g., matrix factorization, PCA, sparse coding and dictionary learning, compressive sensing, etc.)

Why Sparsity?

- A form of regularization. Prevents overfitting
- Fewer parameters to estimate (important when $D \gg N$)
- Can be also seen as automatically performing [feature/variable selection](#)
- Makes it easy to interpret the learned model. Also compact models (need less space to store).
- Predictions can be faster (why?)
- Many natural phenomena are inherently sparse, e.g.,
 - Natural images can be represented as sparse combination of basis images



A Simple Bayesian Approach to Sparsity

- Recall our usual Gaussian prior on \mathbf{w} for linear regression

$$p(\mathbf{w}|\alpha) = \prod_{d=1}^D p(w_d|\alpha) = \prod_{d=1}^D \left(\frac{\alpha}{2\pi}\right)^{1/2} \exp\left(-\frac{\alpha}{2} w_d^2\right)$$

- Each component of \mathbf{w} is a **zero-mean Gaussian** $p(w_d|\alpha) = \mathcal{N}(w_d|0, \alpha^{-1})$
- Same inverse variance (precision) hyperparam α on each entry of \mathbf{w} . Can't impose sparsity on \mathbf{w}
- Let's have a separate inverse variance α_d for each component of \mathbf{w}

$$p(\mathbf{w}|\alpha) = \prod_{d=1}^D p(w_d|\alpha_d) = \prod_{d=1}^D \left(\frac{\alpha_d}{2\pi}\right)^{1/2} \exp\left(-\frac{\alpha_d}{2} w_d^2\right)$$

- We now have D hyperparameters $\alpha = [\alpha_1, \dots, \alpha_D]$ individually controlling the variance of each component w_d of \mathbf{w} . Can infer these using MLE-II, EM, etc (as we've already seen).
 - Note/Caveat: This will not usually yield an exactly sparse solution (i.e., won't get $w_d = 0$)

A Hierarchical Prior to Induce Sparsity

- Our new hierarchical prior on \mathbf{w}

$$p(\mathbf{w}|\boldsymbol{\alpha}) = \prod_{d=1}^D p(w_d|\alpha_d) = \prod_{d=1}^D \left(\frac{\alpha_d}{2\pi}\right)^{1/2} \exp\left(-\frac{\alpha_d}{2} w_d^2\right)$$

- Let's assume a **gamma prior** on each α_d : $p(\alpha_d) \propto \alpha_d^{a-1} \exp^{-\alpha_d/b}$
- The marginal prior on each weight w_d after integrating out α_d

$$p(w_d) = \int p(w_d|\alpha_d)p(\alpha_d)d\alpha_d \quad (\text{will be a Student-t distribution})$$

- Student-t distribution is given by

$$f(t) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi} \Gamma(\frac{\nu}{2})} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

- Akin to penalizing $\sum_{d=1}^D \log |w_d|$. Leads to sparse solutions for \mathbf{w}
- This is often called “Automatic Relevance Determination” (ARD)

Another Sparsity Promoting Prior

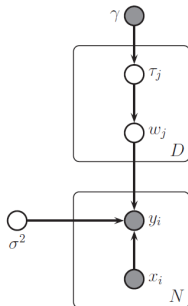
- A **Laplace prior** on weights can also promote sparsity (akin to putting an ℓ_1 regularization)

$$\text{Lap}(w_j|0, 1/\gamma) = \frac{\gamma}{2} e^{-\gamma|w_j|}$$

- Such a prior (and also the previous one) is called **Gaussian scale mixture** prior

$$\text{Lap}(w_j|0, 1/\gamma) = \frac{\gamma}{2} e^{-\gamma|w_j|} = \int \mathcal{N}(w_j|0, \tau_j^2) \text{Ga}(\tau_j^2|1, \frac{\gamma^2}{2}) d\tau_j^2$$

- Basically putting a gamma prior on the variance (a.k.a. “scale”) and integrating it out



An EM algorithm for this model..

- The joint probability distribution of data and all the unknowns

$$p(\mathbf{y}, \mathbf{w}, \boldsymbol{\tau}, \sigma^2 | \mathbf{X}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N) \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{D}_\tau) \propto (\sigma^2)^{-N/2} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2\right) |\mathbf{D}_\tau|^{-\frac{1}{2}} \\ \text{IG}(\sigma^2 | a_\sigma, b_\sigma) \left[\prod_j \text{Ga}(\tau_j^2 | 1, \gamma^2/2) \right] \exp\left(-\frac{1}{2} \mathbf{w}^T \mathbf{D}_\tau \mathbf{w}\right) (\sigma^2)^{-(a_\sigma+1)} \\ \exp(-b_\sigma/\sigma^2) \prod_i \exp(-\frac{\gamma^2}{2} \tau_j^2)$$

- Can also optimize the above using standard optimization methods but using EM has advantages
 - Can use other types of priors on τ_j (leads to other types of sparse priors on \mathbf{w})
 - Easily extend to finding full posterior on \mathbf{w}
 - The technique can be used in a wide variety of models that have ℓ_1 regularization
- Treating τ_j 's as latent variables and \mathbf{w} as parameter, the M step objective is

$$\ell_c(\mathbf{w}) = -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 - \frac{1}{2} \mathbf{w}^\top \mathbf{D}_\tau \mathbf{w} + \text{const}$$

.. easy to optimize (similar to MAP estimation for linear regression)..

An EM algorithm for this model (contd)..

- In the E step, we infer τ_j , whose posterior and the expectation is

$$p(1/\tau_j^2 | \mathbf{w}, \mathcal{D}) = \text{InverseGaussian} \left(\sqrt{\frac{\gamma^2}{w_j^2}}, \gamma^2 \right) \quad \mathbb{E} \left[\frac{1}{\tau_j^2} | w_j \right] = \frac{\gamma}{|w_j|}$$

$$\text{Let } \bar{\Lambda} = \text{diag}(\mathbb{E}[1/\tau_1^2], \dots, \mathbb{E}[1/\tau_D^2])$$

- Note: If we also treat σ^2 as unknown, we can infer it in the E step as

$$p(\sigma^2 | \mathcal{D}, \mathbf{w}) = \text{IG}(a_\sigma + (N)/2, b_\sigma + \frac{1}{2}(\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T(\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})) = \text{IG}(a_N, b_N) \quad \mathbb{E}[1/\sigma^2] = \frac{a_N}{b_N} \triangleq \bar{\omega}$$

- In the M step, we will solve the following problem

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\text{argmax}} -\frac{1}{2}\bar{\omega} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 - \frac{1}{2}\mathbf{w}^T \bar{\Lambda} \mathbf{w}$$

“Hard” Sparsity

- Earlier approaches we saw had a “soft” type of sparsity (α_d/τ_d tells us how important feature d is)
- Consider the basic linear model $\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon$. Denote the data (\mathbf{X}, \mathbf{y}) by \mathcal{D}
- Total D features. Suppose we have indicator variables γ_d , $d = 1, \dots, D$
- $\gamma_d = 1$ if feature d is relevant (w_d is nonzero); $\gamma_d = 0$ otherwise. This is a notion of **hard sparsity**
- The goal is to learn $\gamma = \{\gamma_1, \dots, \gamma_D\}$ (usually along with \mathbf{w})
- If we only want γ , a natural approach would be to do the following MAP estimation

$$\hat{\gamma} = \arg \max_{\gamma} p(\gamma|\mathcal{D}) = \arg \max_{\gamma} p(\gamma)p(\mathcal{D}|\gamma)$$

where $p(\gamma)$ is the prior on γ and $p(\mathcal{D}|\gamma) = \int p(\mathcal{D}|\gamma, \mathbf{w})p(\mathbf{w})d\mathbf{w}$ is the marginal likelihood

- A combinatorial search problem. A naïve approach to infer $\hat{\gamma}$ would require searching over a set of size 2^D . But there are more efficient ways (e.g., greedy search)
- But first we need to write down the **marginal posterior over γ** , $p(\gamma|\mathcal{D}) \propto p(\gamma)p(\mathcal{D}|\gamma)$

The Prior

- The marginal posterior $p(\gamma|\mathcal{D}) \propto p(\gamma)p(\mathcal{D}|\gamma)$
- Let's consider the following **prior** on the binary vector $\gamma = \{\gamma_1, \dots, \gamma_D\}$

$$p(\gamma) = \prod_{j=1}^D \text{Ber}(\gamma_j|\pi_0) = \pi_0^{||\gamma||_0} (1 - \pi_0)^{D - ||\gamma||_0}$$

- π_0 is the probability that a feature is relevant (corresponding γ_d is nonzero)
- $||\gamma||_0 = \sum_{d=1}^D \gamma_d$ is the number of nonzeros in γ (the ℓ_0 norm of γ)
- Note the form of the log-prior

$$\begin{aligned} \log p(\gamma|\pi_0) &= ||\gamma||_0 \log \pi_0 + (D - ||\gamma||_0) \log(1 - \pi_0) \\ &= ||\gamma||_0 (\log \pi_0 - \log(1 - \pi_0)) + \text{const} \\ &= -\lambda ||\gamma||_0 + \text{const} \end{aligned}$$

where $\lambda \triangleq \log \frac{1-\pi_0}{\pi_0}$ controls the sparsity of the model.

The (Marginal) Likelihood

- Let's now look at the marginal **likelihood**

$$p(\mathcal{D}|\gamma) = p(\mathbf{y}|\mathbf{X}, \gamma) = \int \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \gamma) p(\mathbf{w}|\gamma, \sigma^2) p(\sigma^2) d\mathbf{w} d\sigma^2$$

- Given γ , the prior $p(\mathbf{w}|\gamma, \sigma^2)$ on \mathbf{w}

$$p(w_j|\sigma^2, \gamma_j) = \begin{cases} \delta_0(w_j) & \text{if } \gamma_j = 0 \\ \mathcal{N}(w_j|0, \sigma^2 \sigma_w^2) & \text{if } \gamma_j = 1 \end{cases}$$

where w_j is the j -th element of \mathbf{w}

- This is also called a **spike-and-slab prior**
- Ignoring element of \mathbf{w} that are zero, we get following marginal likelihood

$$p(\mathcal{D}|\gamma) = \int \int \mathcal{N}(\mathbf{y}|\mathbf{X}_\gamma \mathbf{w}_\gamma, \sigma^2 \mathbf{I}_N) \mathcal{N}(\mathbf{w}_\gamma|\mathbf{0}_{D_\gamma}, \sigma^2 \sigma_w^2 \mathbf{I}_{D_\gamma}) p(\sigma^2) d\mathbf{w}_\gamma d\sigma^2$$

where D_γ is the number of nonzeros in γ and \mathbf{X}_γ consists of the corresponding columns of \mathbf{X}

The (Marginal) Likelihood

- Still have the noise variance hyperparameter σ^2
- If σ^2 is known then

$$p(\mathcal{D}|\gamma, \sigma^2) = \int \mathcal{N}(\mathbf{y}|\mathbf{X}_\gamma \mathbf{w}_\gamma, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w}_\gamma|\mathbf{0}, \sigma^2 \Sigma_\gamma) d\mathbf{w}_\gamma = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_\gamma)$$
$$\mathbf{C}_\gamma \triangleq \sigma^2 \mathbf{X}_\gamma \Sigma_\gamma \mathbf{X}_\gamma^T + \sigma^2 \mathbf{I}_N$$

- If σ^2 is unknown, we can put an inverse gamma prior $p(\sigma^2) = \text{IG}(\sigma^2|a_\sigma, b_\sigma)$ and integrate out σ^2

$$p(\mathcal{D}|\gamma) = \int \int p(\mathbf{y}|\gamma, \mathbf{w}_\gamma, \sigma^2) p(\mathbf{w}_\gamma|\gamma, \sigma^2) p(\sigma^2) d\mathbf{w}_\gamma d\sigma^2$$
$$\propto |\mathbf{X}_\gamma^T \mathbf{X}_\gamma + \Sigma_\gamma^{-1}|^{-\frac{1}{2}} |\Sigma_\gamma|^{-\frac{1}{2}} (2b_\sigma + S(\gamma))^{-(2a_\sigma + N - 1)/2}$$

$S(\gamma)$ denotes the residual sum-of-squares $\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}_\gamma (\mathbf{X}_\gamma^T \mathbf{X}_\gamma + \Sigma_\gamma^{-1})^{-1} \mathbf{X}_\gamma^T \mathbf{y}$

Solving for γ

- Now we have an expression (exact/approximate) for $(\log)p(\mathcal{D}|\gamma)$
- This combined with $\log p(\gamma) \propto -\lambda\|\gamma\|_0$ gives the posterior over γ
- Now the goal is to find γ
- One way is to do it by greedily choosing the next feature to add
- Some heuristics
 - Single best replacement (start with $\gamma = 0$ and consider **all one-bit flips**)
 - Matching pursuit (for ordinary least squares problems)

$$j^* = \arg \min_{j \notin \gamma_t} \min_{\beta} \|\mathbf{y} - \mathbf{X}\mathbf{w}_t - \beta \mathbf{x}_{:,j}\|^2$$

$$\beta = \mathbf{x}_{:,j}^T \mathbf{r}_t / \|\mathbf{x}_{:,j}\|^2$$

$$j^* = \arg \max \mathbf{x}_{:,j}^T \mathbf{r}_t$$

where $\mathbf{r}_t = \|\mathbf{y} - \mathbf{X}\mathbf{w}_t\|$ is the current residual. Thus we choose the feature (a column in the $N \times D$ feature matrix \mathbf{X}) that has the maximum correlation with current residual

- Do also do fully Bayesian inference using MCMC or VB