

$$5.1 \quad \text{To solve} \quad \left. \begin{aligned} & \min \frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2 \\ & \text{s.t.} \quad \sum_{j=1}^n |w_j|^q \leq \eta \end{aligned} \right\} \quad \textcircled{1}$$

Use Lagrange multipliers,

$$L = \frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \alpha \left(\sum_{j=1}^n |w_j|^q - \eta \right)$$

$$w^* = \min_w \left[\max_{\alpha: \alpha \geq 0} \frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \alpha \left(\sum_{j=1}^n |w_j|^q - \eta \right) \right]$$

$$= \min_w \left[\frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \alpha^* \sum_{j=1}^n |w_j|^q - \alpha^* \eta \right]$$

$$= \min_w \left[\frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \alpha^* \sum_{j=1}^n |w_j|^q \right]$$

It means when $\lambda = 2\alpha^*$, the minimization of the regularized error function is equivalent to ①.

$$5.2 \quad \max_{\theta} \prod_{i=1}^n P(x_i|\theta) P(\theta)$$

$$\Rightarrow \max_{\theta} \left[\log \prod_{i=1}^n P(x_i|\theta) + \log P(\theta) \right]$$

$$\Rightarrow \max_{\theta} \left[\log \prod_{i=1}^n P(x_i|\theta) + \log \frac{\lambda}{2} - \lambda |\theta| \right]$$

$$\Rightarrow \max_{\theta} \left[\log \prod_{i=1}^n P(x_i|\theta) - \lambda |\theta| \right]$$

$$\begin{aligned}
 5.3(a). \text{MSE}(\hat{x}) &= E[(\hat{x} - x)^2] \\
 &= E[(\hat{x} - x)^T(\hat{x} - x)] \\
 &= E[\hat{x}^T \hat{x}] - 2E[\hat{x}^T x] + E[x^T x]
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(\hat{x}) &= E[(\hat{x} - E[\hat{x}])^2] \\
 &= E[(\hat{x} - E[\hat{x}])^T(\hat{x} - E[\hat{x}])] \\
 &= E[\hat{x}^T \hat{x} - \hat{x} E[\hat{x}] - E[\hat{x}] \hat{x} + E[\hat{x}]^T E[\hat{x}]] \\
 &= E[\hat{x}^T \hat{x} + E[\hat{x}]^T E[\hat{x}] - 2E[\hat{x}]^T \hat{x}]
 \end{aligned}$$

$$\begin{aligned}
 [\text{Bias}(\hat{x})]^2 &= (E[\hat{x}] - x)^T(E[\hat{x}] - x) \\
 &= E[\hat{x}]^T E[\hat{x}] - E[\hat{x}]^T x - x^T E[\hat{x}] + x^T x
 \end{aligned}$$

$$\text{Var}(\hat{x}) + (\text{Bias}(\hat{x}))^2 = E(\hat{x}^T \hat{x}) - 2E[\hat{x}]^T x + x^T x$$

Also, we have $E[x] = x$ and $E[\hat{x}^T x] = E[\hat{x}]^T x$

$$\therefore \text{Var}(\hat{x}) + (\text{Bias}(\hat{x}))^2 = \text{MSE}(\hat{x})$$

$$\begin{aligned}
 (b) \text{E}[(\hat{x} - x - N)^2] &= E[(\hat{x} - x - N)^T(\hat{x} - x - N)] \\
 &= E[(\hat{x} - x)^T(\hat{x} - x) - 2(\hat{x} - x)^T N + N^T N] \\
 &= \text{MSE}(\hat{x}) - 2E[(\hat{x} - x)^T E(N)] + E[N^T N] \\
 &= \text{MSE}(\hat{x}) + S^2
 \end{aligned}$$

5.4 (a) $\min(k, |X| - k)$

Since for hypothesis set H_k^x , there are exactly k elements in X are labeled as 1 and $(|X| - k)$ elements are labeled as 0. So H can scatter any $\min(k, |X| - k)$ points.

For ~~any points~~ ^{If there are} $(k+1)$ points labeled 1 or $(|X| - k + 1)$ points labeled as 0, H_k^x cannot scatter them.

It means $VC(H_k^x) < k+1$ and $VC(H_k^x) < |X| - k + 1$

So $VC(H_k^x) = \min(k, |X| - k)$

(b). $\min(2k+1, |X|)$

For $H_{\leq k}^x$, there are at most k elements of X to make $h(x)=1$ or $h(x)=0$. It means when we give a set of points in X , and there always exists ^m ~~at least~~ points labeled 1 or 0, where $m \leq k$, then we can find a hypothesis from $H_{\leq k}^x$ to scatter them.

So for any points less than $(2k+1)$, $H_{\leq k}^x$ can scatter them.

For $(2k+2)$ points with $(k+1)$ labeled 1 and $(k+1)$ labeled 0, $H_{\leq k}^x$ cannot scatter them. $VC(H_{\leq k}^x) < (2k+2)$

Also, we should consider $|X|$, so $VC(H_{\leq k}^x) = \min(2k+1, |X|)$.

Large Scale Linear Programming Decoding via the Alternating Direction Method of Multipliers

The keynote speech of **Stark Draper**

Professor shared his application of the ADMM to solve the linear programming relaxation of maximum likelihood decoding for error-correction codes. ADMM, which is a topic I am interested about. I am curious about the algorithm but the area of error-correction codes is unfamiliar to me. So I was lost in the most important part of content. But the background knowledge part gave me some intuitions about the ADMM.

The objective function is separable in x and y . The dual update requires solving a proximity function in x and y at the same time; the ADMM technique allows this problem to be solved approximately by first solving for x with y fixed, and then solving for y with x fixed. Rather than iterate until convergence, the algorithm proceeds directly to updating the dual variable and then repeating the process in an efficient and parallelizable manner. Because of this approximation, the algorithm is distinct from the pure augmented Lagrangian method.

To summarize, ADMM is an algorithm that solves convex optimization problems by breaking them into smaller pieces, each of which are then easier to handle. It has recently found wide application in a number of areas.