

$$1/ \quad J(x) = \|Ax - y\|_2^2 + \lambda \|x\|_2^2 \quad A \in \mathbb{R}^{N \times D}, x \in \mathbb{R}^D$$

$$= (Ax - y)^T (Ax - y) + \lambda x^T x$$

$$g(x) = \frac{\partial J}{\partial x} = \frac{\partial (x^T A^T A x - 2x^T A^T y + y^T y + \lambda x^T x)}{\partial x}$$

$$= 2A^T A x - 2A^T y + 2\lambda x$$

For exact line search, the optimal step size for  $n$ -th step:

$$\mu_n^* = \underset{\mu \geq 0}{\operatorname{argmin}} J(x_{n-1} - \mu g(x_{n-1})) \quad (*)$$

Consider  $J(x_{n-1} - \mu g(x_{n-1}))$

$$= \|A(x_{n-1} - \mu g(x_{n-1})) - y\|_2^2 + \lambda \|x_{n-1} - \mu g(x_{n-1})\|_2^2$$

$$\bullet \|A(x_{n-1} - \mu g(x_{n-1})) - y\|_2^2$$

$$= \|Ax_{n-1} - \mu Ag(x_{n-1}) - y\|_2^2$$

$$= (Ax_{n-1} - \mu Ag(x_{n-1}) - y)^T (Ax_{n-1} - \mu Ag(x_{n-1}) - y)$$

$$= x_{n-1}^T A^T A x_{n-1} - x_{n-1}^T A^T \mu Ag(x_{n-1}) - x_{n-1}^T A^T y$$

$$- \mu g^T(x_{n-1}) A^T A x_{n-1} + \mu g^T(x_{n-1}) A^T \mu Ag(x_{n-1})$$

$$+ \mu g^T(x_{n-1}) A^T y - y^T A x_{n-1} + y^T \mu Ag(x_{n-1}) + y^T y$$

By denoting  $D = A^T A$ , we have:

$$\|A(x_{n-1} - \mu g(x_{n-1})) - y\|_2^2$$

$$= x_{n-1}^T D x_{n-1} - \mu x_{n-1}^T D g(x_{n-1}) - x_{n-1}^T A^T y$$

$$- \mu g^T(x_{n-1}) D x_{n-1} + \mu^2 g^T(x_{n-1}) D g(x_{n-1})$$

$$+ \mu g^T(x_{n-1}) A^T y - y^T A x_{n-1} + \mu y^T A g(x_{n-1}) + y^T y$$

$$= \mu^2 g^T(u_{n-1}) Dg(u_{n-1}) + \mu [-u_{n-1}^T Dg(u_{n-1}) - g^T(u_{n-1}) D u_{n-1} + g^T(u_{n-1}) A^T y + y^T A g(u_{n-1})] + C_1$$

with  $C_1 = u_{n-1}^T D u_{n-1} - u_{n-1}^T A^T y - y^T A u_{n-1} + g^T y$

(The reason why do not write  $C$  explicitly is we will get the gradient with respect to  $\mu$ , therefore we only concern terms having  $\mu$ ).

$$\begin{aligned} & \lambda \|u_{n-1} - \mu g(u_{n-1})\|_2^2 \\ &= \lambda (u_{n-1} - \mu g(u_{n-1}))^T (u_{n-1} - \mu g(u_{n-1})) \\ &= \lambda (u_{n-1}^T u_{n-1} - u_{n-1}^T \mu g(u_{n-1}) - \mu g^T(u_{n-1}) u_{n-1} + \mu^2 g^T(u_{n-1}) g(u_{n-1})) \\ &= \mu^2 g^T(u_{n-1}) g(u_{n-1}) \lambda - 2\mu u_{n-1}^T g(u_{n-1}) \lambda + C_2 \end{aligned}$$

with  $C_2 = \lambda u_{n-1}^T u_{n-1}$

The objective function in (\*) become:

$$\begin{aligned} J(\mu) &= \mu^2 [g^T(u_{n-1}) Dg(u_{n-1}) + \lambda g^T(u_{n-1}) g(u_{n-1})] \\ &\quad + \mu [-2u_{n-1}^T Dg(u_{n-1}) + 2g^T(u_{n-1}) A^T y - 2\lambda u_{n-1}^T g(u_{n-1})] + C_1 + C_2 \end{aligned}$$

$J(\mu)$  is quadratic function with  $\mu$ , therefore we are able to solve (\*) by stationary condition, then check the found  $\mu^*$  satisfy greater than equal 0 or not.

$$\begin{aligned} \frac{\partial J(\mu)}{\partial \mu} &= 2\mu [g^T(u_{n-1}) Dg(u_{n-1}) + \lambda g^T(u_{n-1}) g(u_{n-1})] \\ &\quad - 2[u_{n-1}^T Dg(u_{n-1}) - g^T(u_{n-1}) A^T y + \lambda u_{n-1}^T g(u_{n-1})] \end{aligned}$$

By setting  $\frac{\partial J(\mu)}{\partial \mu} = 0$ , we have:

$$\begin{aligned}\mu_n &= \frac{w_{n-1}^T D g(w_{n-1}) - g^T(w_{n-1}) A^T y + \lambda w_{n-1}^T g(w_{n-1})}{g^T(w_{n-1}) D g(w_{n-1}) + \lambda g^T(w_{n-1}) g(w_{n-1})} \\ &= \frac{w_{n-1}^T (D + \lambda I) g(w_{n-1}) - g^T(w_{n-1}) A^T y}{g^T(w_{n-1}) (D + \lambda I) g(w_{n-1})}.\end{aligned}$$

In summary, in the step  $n-1$  of Gradient Descent, we got  $w_{n-1} \in \mathbb{R}^{D \times 1}$ ,  $g(w_{n-1}) \in \mathbb{R}^{D \times 1}$  by plugging  $w_{n-1}$  into the first formula appearing in solution. Then we calculate the optimal step size  $\mu_n^*$  by:

$$\mu_n^* = \frac{w_{n-1}^T (D + \lambda I) g(w_{n-1}) - g^T(w_{n-1}) A^T y}{g^T(w_{n-1}) (D + \lambda I) g(w_{n-1})}$$

Then update the next parameter  $w_n$  by:

$$w_n = w_{n-1} - \mu_n^* g(w_{n-1})$$

→ The objective function  $J(x) = \|Ax - y\|_2^2 + \lambda \|x\|_2^2$  is also called ridge regression by putting additional term  $\lambda \|x\|_2^2$  into estimate mean square error ( $\|Ax - y\|_2^2$ ).

→ By including  $\lambda$ , we are able to handle the linearly multicollinearity dependence in data which make  $A^T A$  is not invertible ( $x = (A^T A)^{-1} A^T y$ , also known as closed form solution). When we want to find the exact solution in 1 shot.

Also, the characteristic of  $L_2$ -norm is uniformity which restrict all element  $x_i$  to low value and therefore can help reducing overfitting.

End Q.1



$$\alpha / \quad s(x) = [f(x) - \hat{f}(x)]^2$$

Data model:

$$y(x) = f(x) + \epsilon \quad (f \text{ is unknown, deterministic})$$

Trained function  $\hat{f}(\cdot)$  that estimates  $f(x)$

$$\rightarrow \text{Bias}_S(\hat{f}(x)) = f(x) - E_S[\hat{f}(x)]$$

$$\rightarrow \text{Var}_S(\hat{f}(x)) = E_S[(\hat{f}(x) - E_S[\hat{f}(x)])^2]$$

$$\rightarrow \text{MSE}_{S, \epsilon}(\hat{f}_S(x)) = E_{S, \epsilon}[s(x)]$$

$$= E_{S, \epsilon}[(y(x) - \hat{f}(x))^2]$$

We need to prove

$$\text{MSE}_{S, \epsilon}(\hat{f}_S(x)) = \text{Bias}_S^2(\hat{f}(x)) + \text{Var}_S(\hat{f}(x))$$

We have:

$$\text{MSE}_{S, \epsilon} = E[(y - \hat{f})^2] \stackrel{y=f+\epsilon}{=} E[(f + \epsilon - \hat{f})^2]$$

$$= E[(f - \hat{f})^2 + \epsilon^2 + 2\epsilon(f - \hat{f})]$$

$$= E\{(f - \hat{f})^2\} + E\{\epsilon^2\} + 2E\{(f - \hat{f})\epsilon\}$$

linearity  
of Expectation  
operator

$$\text{We have: } \text{var}(\epsilon) = E(\epsilon^2) - E(\epsilon)^2$$

$$\sigma_\epsilon^2 = E(\epsilon^2) - 0 \Rightarrow E\{\epsilon^2\} = \sigma_\epsilon^2$$

$$E\{(f - \hat{f})\epsilon\} = E\{f - \hat{f}\} E\{\epsilon\} = E\{f - \hat{f}\} 0 = 0$$

independence

Therefore:

$$MSE_{S, \theta} = E\{(f - \hat{f})^2\} + \sigma_e^2$$

Consider  $E\{(f - \hat{f})^2\}$ , we have:

$$\begin{aligned} E\{(f - \hat{f})^2\} &= E\{f^2 - 2f\hat{f} + \hat{f}^2\} \\ &= E\{f^2\} - 2E\{f\hat{f}\} + E\{\hat{f}^2\} \end{aligned}$$

$$\begin{aligned} \text{Because } f \text{ is deterministic } \Rightarrow \quad & E\{f^2\} = f^2 \\ & E\{f\hat{f}\} = fE\{\hat{f}\} \end{aligned}$$

$$\Rightarrow E\{(f - \hat{f})^2\} = f^2 - 2fE\{\hat{f}\} + E\{\hat{f}^2\}$$

$$\begin{aligned} &= f^2 + E\{\hat{f}^2\} - 2fE\{\hat{f}\} + E\{\hat{f}\}^2 - E\{\hat{f}\}^2 \\ &\quad - 2E\{\hat{f}\}^2 + 2E\{\hat{f}\}^2 \end{aligned}$$

$$= f^2 + E[\hat{f}^2 + E\{\hat{f}\}^2 - 2E\{\hat{f}\}^2] - 2fE\{\hat{f}\} + E\{\hat{f}\}^2$$

$$= [f^2 - 2fE\{\hat{f}\} + E\{\hat{f}\}^2] + E[\hat{f}^2 + E\{\hat{f}\}^2 - 2E\{\hat{f}\}^2]$$

$$= (f - E\{\hat{f}\})^2 + E\{(f - E\{\hat{f}\})^2\}$$

$$= \text{Bias}^2(\hat{f}) + \text{Var}(\hat{f})$$

End Q2

3/

- The term "The curse of dimensionality" in non-parametric learning refers to challenges in high-dimensional space (i.e.,  $d$  grows fast)
- One of parts was introduced is curse of sample size. It means that as the number of dimensions  $d$  increases, the fraction of volume considered as neighborhoods would be increase exponentially
- Several challenges :
  - + computational complexity :  $d$  grows fast  $\rightarrow$  cost for calculating distance between data point increase exponentially.
  - + Increased sparsity of data :  $d \uparrow$  then sparsity  $\uparrow$  which leads unreliable density estimation, hard to determine the relationships between variables accurately.
  - + Overfitting and generalization :  $d \uparrow$ ,  $P_n(\text{overfitting}) \uparrow$ , therefore need to provide more training data  $\rightarrow$  poor generalization on unseen data.
- Some techniques to mitigate "Curse of Dimensionality"
  - + Dimensionality reduction (LDA, PCA) }
  - + Feature selection }
  - + Regularization ( $L1 \rightarrow$  make model sparsity,  $L2$ ) }

#### 4/ Parametric model

$$y = b(x)^T \kappa + \varepsilon, \quad \kappa \in \mathbb{R}^D$$

$\{(y_i, x_i)\}_{i=1}^N$  : training data.

Prior :

$$\kappa \sim N(m_0, C_0), \quad m_0 \in \mathbb{R}^D, \quad C_0 \in \mathbb{R}^{D \times D}$$

$$\text{or } f_{\#}(\kappa) = N(m_0, C_0)$$

MAP try to solve:

$$\max_{\kappa} f_{\#}(\kappa | y, X)$$

By Bayes rule, we can rewrite  $f_{\#}(\kappa | y, X)$  as:

$$f_{\#}(\kappa | y, X) = \frac{f_{\#}(y, X | \kappa) f_{\#}(\kappa)}{f_{\#}(y, X)}$$

$$= \frac{f_{\#}(y | X, \kappa) f_{\#}(X | \kappa) f_{\#}(\kappa)}{f_{\#}(y, X)}$$

$$= \frac{f_{\#}(y | X, \kappa) f_{\#}(X) f_{\#}(\kappa)}{f_{\#}(y, X)}$$

$$= f_{\#}(y | X, \kappa) f_{\#}(\kappa) c$$

$$\text{with } c = \frac{f_{\#}(X)}{f_{\#}(y, X)} > 0$$

therefore  $\max_{\kappa} f_{\#}(\kappa | y, X)$  is written equivalently as:

$$\max_{\kappa} f_{\#}(y | X, \kappa) f_{\#}(\kappa)$$



$$\equiv \max_{\mathbf{x}} \ln(f_{**}(y | X, \mathbf{x})) + \ln(f_{**}(\mathbf{x})) \quad (**)$$

By MLE, we found:

$$\mathbf{H} = [b(x_1), \dots, b(x_n)]$$

$$f_{**}(y | X, \mathbf{x}) = N(\mathbf{H}^T \mathbf{x}, V \mathbf{I}_N)$$

$$c_1 > 0$$

$$= c_1 \exp \left\{ -\frac{1}{2} (y - \mathbf{H}^T \mathbf{x})^T (V \mathbf{I}_N)^{-1} (y - \mathbf{H}^T \mathbf{x}) \right\}$$

By prior distribution:

$$f_{**}(\mathbf{x}) = N(\mathbf{m}_0, \mathbf{C}_0) = c_2 \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mathbf{m}_0)^T \mathbf{C}_0^{-1} (\mathbf{x} - \mathbf{m}_0) \right\}$$

$$c_2 \geq 0$$

$$\Rightarrow \ln(f_{**}(y | X, \mathbf{x})) = \ln(c_1) - \frac{1}{2V} \|\mathbf{y} - \mathbf{H}^T \mathbf{x}\|_2^2$$

$$\Rightarrow \ln(f_{**}(\mathbf{x})) = \ln(c_2) - \frac{1}{2} (\mathbf{x} - \mathbf{m}_0)^T \mathbf{C}_0^{-1} (\mathbf{x} - \mathbf{m}_0)$$

(\*\*) become:

$$\max_{\mathbf{x}} \ln(c_1) + \ln(c_2) - \frac{1}{2V} \|\mathbf{y} - \mathbf{H}^T \mathbf{x}\|_2^2 - \frac{1}{2} (\mathbf{x} - \mathbf{m}_0)^T \mathbf{C}_0^{-1} (\mathbf{x} - \mathbf{m}_0)$$

$$\equiv \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{H}^T \mathbf{x}\|_2^2 + (\mathbf{x} - \mathbf{m}_0)^T (V \mathbf{C}_0^{-1}) (\mathbf{x} - \mathbf{m}_0)$$

$\|\mathbf{y} - \mathbf{H}^T \mathbf{x}\|_2^2$  is quadratic

$(\mathbf{x} - \mathbf{m}_0)^T (V \mathbf{C}_0^{-1}) (\mathbf{x} - \mathbf{m}_0)$  is quadratic

$\Rightarrow$  the objective is convex

By applying Stationary condition :

$$\frac{\partial T(x)}{\partial x} = \frac{\partial (y^T y + x^T H H^T x - 2x^T H y + x^T (V C_0^{-1}) x + m_0^T (V C_0^{-1}) m_0)}{\partial x}$$

$$= 2H H^T x - 2H y + 2(V C_0^{-1}) x - 2(V C_0^{-1}) m_0$$

$\stackrel{\text{set}}{=} 0$

$$\Rightarrow \boxed{x_{\text{MAP}} = (H H^T + V C_0^{-1})^{-1} (H y + V C_0^{-1} m_0)}$$

↑ map solution

The solution in MLE

$$x_{\text{MLE}} = (H H^T)^{-1} H y$$

We can see that the MAP work when we have both knowledge about prior distribution and given dataset.

In case of MLE, we only know training data  $\{(y_i, x_i)\}_{i=1}^N$  while in MAP, we are given prior distribution of  $x$ .

→ MAP is generalized version of MLE.

→ If we don't know prior knowledge, we boil down to MLE by considering  $C_0 = a I$  with  $a \rightarrow \infty$

$$x_{\text{MAP}} \xrightarrow{C_0^{-1} \rightarrow 0} (H H^T)^{-1} H y = x_{\text{MLE}}$$

→ If we do not have dataset then we turn into prior distribution maximization

$$H \rightarrow 0 \Rightarrow x_{\text{MAP}} = (0 + V C_0^{-1})^{-1} (0 + C_0^{-1} m_0 V) \\ = (V C_0^{-1})^{-1} (V C_0^{-1}) m_0 = m_0$$

→ Another view of MAP is least square + regularization

$$\min_x \|y - H^T x\| + \underbrace{(x - m_0)^T (V C_0^{-1}) (x - m_0)}$$

↓  
least square error

↓  
penalized term  
by prior knowledge.

End Q4

5/ Binary classification with logistic regression and training data  
 $\{(t_n, x_n)\}_{n=1}^N$

Cross-entropy loss :

$$\begin{aligned} L(w) &= - \sum_{n=1}^N (\ln(y_n) t_n + \ln(1-y_n)(1-t_n)) \\ &= - \sum_{n=1}^N (\ln(\sigma(w^T b(x_n))) t_n + \ln(1-\sigma(w^T b(x_n)))(1-t_n)) \end{aligned}$$

The gradient of  $L(w)$  :

$$\begin{aligned} \nabla_w L(w) &= \nabla_w - \sum_{n=1}^N \left[ \ln(\sigma(w^T b(x_n))) t_n + \ln(1-\sigma(w^T b(x_n)))(1-t_n) \right] \\ &= - \sum_{n=1}^N \left( \nabla_w \ln(\sigma(w^T b(x_n))) t_n + \nabla_w \ln(1-\sigma(w^T b(x_n)))(1-t_n) \right) \end{aligned}$$

Consider  $\nabla_w \ln(\sigma(w^T b(x_n)))$  :

$$= \frac{1}{\sigma(w^T b(x_n))} \nabla_w \sigma(w^T b(x_n))$$

$$= \frac{1}{\sigma(w^T b(x_n))} \sigma'(w^T b(x_n)) (1 - \sigma(w^T b(x_n))) b(x_n)$$

$$= (1 - \sigma(w^T b(x_n))) b(x_n)$$

Consider  $\nabla_w \ln(1 - \sigma(w^T b(x_n)))$



$$= \frac{1}{1 - \sigma'(\mathbf{x}^T \mathbf{b}(x_n))} \nabla_{\mathbf{x}} (1 - \sigma'(\mathbf{x}^T \mathbf{b}(x_n)))$$

$$= \frac{1}{1 - \sigma'(\mathbf{x}^T \mathbf{b}(x_n))} \nabla_{\mathbf{x}} (-\sigma'(\mathbf{x}^T \mathbf{b}(x_n)))$$

$$= \frac{-1}{1 - \sigma'(\mathbf{x}^T \mathbf{b}(x_n))} \sigma'(\mathbf{x}^T \mathbf{b}(x_n)) (1 - \sigma'(\mathbf{x}^T \mathbf{b}(x_n))) \mathbf{b}(x_n)$$

$$= -\sigma'(\mathbf{x}^T \mathbf{b}(x_n)) \mathbf{b}(x_n)$$

Therefore :

$$\begin{aligned} \nabla_{\mathbf{x}} L(\mathbf{x}) &= - \sum_{n=1}^N (1 - \sigma'(\mathbf{x}^T \mathbf{b}(x_n)) \mathbf{b}(x_n) t_n \\ &\quad - \sigma'(\mathbf{x}^T \mathbf{b}(x_n)) \mathbf{b}(x_n) (1 - t_n)) \\ &= - \sum_{n=1}^N \mathbf{b}(x_n) [(1 - \sigma'(\mathbf{x}^T \mathbf{b}(x_n)) t_n \\ &\quad - \sigma'(\mathbf{x}^T \mathbf{b}(x_n)) (1 - t_n)] \\ &= - \sum_{n=1}^N \mathbf{b}(x_n) [t_n - t_n \sigma'(\mathbf{x}^T \mathbf{b}(x_n)) - \sigma'(\mathbf{x}^T \mathbf{b}(x_n)) \\ &\quad + t_n \sigma'(\mathbf{x}^T \mathbf{b}(x_n))] \\ &= - \sum_{n=1}^N \mathbf{b}(x_n) (t_n - \sigma'(\mathbf{x}^T \mathbf{b}(x_n))) \end{aligned}$$

Therefore

$$g(\mathbf{x}) = \nabla_{\mathbf{x}} L(\mathbf{x}) = - \sum_{n=1}^N \mathbf{b}(x_n) (t_n - \sigma'(\mathbf{x}^T \mathbf{b}(x_n)))$$

End Q5

6/

### Adam algorithm

1/ Initialize  $W_0$ ;  $s_0 = 0$ ;  $\Delta_0 = 0$

2/ For  $i = 1, 2, \dots$  until termination criteria holds

3/ Update  $s_i$ :  $[s_i]_j = \gamma [s_{i-1}]_j + (1-\gamma) [g_i]_j^2, \forall j \in [d]$

4/ Bias correction  $s$ :  $s = s_i \frac{1}{1-\gamma^i}$

5/ Update  $V_i = \alpha V_{i-1} + (1-\alpha) g_i$

6/ Bias correction  $V$ :  $V = V_i \frac{1}{1-\alpha^i}$

7/ Update parameter  $W_i$ :  $W_i = W_{i-1} - \mu (\text{diag}(s)^{\frac{1}{2}} + \epsilon I)^{-1} V$

8/ Return  $W_i$

### Potential merits of Adam:

1/ It combines benefits of Momentum SGD and RMSProp

- RMSProp characteristic: adaptively change the learning rate

- Momentum characteristic: accelerate convergence.

2/ It includes bias correction which helps fixing the bias when initializing moment estimates at zero  $\rightarrow$  more accurate estimations, specially in early steps of training.

3/ Handles noisy data and works well with sparse gradients more effectively compared to Momentum SGD and RMSProp

4/ Can get good performance without too much effort in tuning hyperparameters.

End Q6

7/ K-class Softmax Logistic Regression and training data  
 $\{(t_n, x_n)\}_{n=1}^N$

Cross-entropy loss

$$L(W) = -\ln(p(T|H, W))$$

$$= -\ln\left(\prod_{n=1}^N \prod_{k=1}^K [Y(W)]_{k,n}^{[T]_{k,n}}\right)$$

$$= -\sum_{n=1}^N \sum_{k=1}^K [T]_{k,n} \cdot \ln[Y(W)]_{k,n}$$

with  $W \in \mathbb{R}^{D \times K}$

Gradient of  $L(W)$   $\rightarrow G(W)$

$$\mathbb{R}^{D \times K} \ni G(W) = [g_1(W), \dots, g_K(W)]$$

$$= [H a_1(W), \dots, H a_K(W)]$$

$$= H A(W)$$

$$A(W) = [a_1(W), \dots, a_K(W)] = (Y(W) - T)^T$$

$$\Rightarrow G(W) = H (Y(W) - T)^T$$

with  $H = [b(x_1), \dots, b(x_N)]$ ,  $[T]_{:,i} = \{[t_1]_i, \dots, [t_N]_i\}^T$   
 $[Y(W)]_{:,i} = [y_1(W)]_i, \dots, [y_N(W)]_i]^T$

End Q7.

8/

$$\{x_n \in \mathbb{R}^D\}_{n=1}^N$$

L1-PCA formula:

$$\max_{G \in \mathbb{R}^{D \times K}, G^T G = I_K} \sum_{n=1}^N \|G^T x_n\|_1$$

Fixed point iterations algorithm:

1/ Initialize  $B \in \{-1, 1\}^{N \times K}$  arbitrarily

2/ Repeat until termination conditions hold

3/  $U, S, V^T \xleftarrow{\text{thin SVD}} (HB)$

4/  $G \leftarrow UV^T$

5/  $B \leftarrow \text{sign}(H^T G)$

6/ Return  $G$

End Q8

EE 6363. Advanced Machine Learning

Teen Van Nguyen

xdp 076