

1. Let  $u$  be a unit-length vector and  $x \in H$ . Then

$$\|x - xu\|^2 = \|x\|^2 - 2\langle x, u \rangle + u^2$$

For fixed  $x$ , if this is a quadratic expression in  $u$ , that takes its minimum iff  $u = \langle x, u \rangle$

Therefore:

$$f_u(x) = \arg \min_{u \in V} \|x - vu\|$$

Hence:

$$\begin{aligned} \sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|^2 &= \sum_{i=1}^n \|x^{(i)}\|^2 - 2\langle x^{(i)}, f_u(x^{(i)}) \rangle + \|\langle x^{(i)}, u \rangle\|^2 \\ &= \sum_{i=1}^n \|x^{(i)}\|^2 - 2\langle x^{(i)}, \langle x^{(i)}, u \rangle u \rangle + \|\langle x^{(i)}, u \rangle u\|^2 \end{aligned}$$

$$\begin{aligned}
 &= \sum_{i=1}^n \|x^{(i)}\|^2 - 2\langle x^{(i)}, u \rangle^2 + \langle x^{(i)}, u \rangle^2 \\
 &= \sum_{i=1}^n \|x^{(i)}\|^2 - \langle x^{(i)}, u \rangle^2
 \end{aligned}$$

Let's conclude that!

$$\arg \min_{u \in \mathbb{R}^n} \sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|^2 = \arg \max_{u \in \mathbb{R}^n} \sum_{i=1}^n \langle x^{(i)}, u \rangle^2$$

Where the right side is exactly the dot product of the first principal component, i.e. the vector that maximizes the variance of the projections of the data.

L.2

As we know from the lecture notes that

$$J(\theta, \vartheta) := \sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})},$$

where  $Q$  is any set of distributions  $Q_i$ , we have the inequality

$$J(\theta) \geq J(\theta, \vartheta),$$

with respect to

in which equality holds iff

$$Q_i(z^{(i)}) = p(x^{(i)} | z^{(i)}; \theta).$$

Using J we can write  $\theta^{(t+1)}$  as:

$$\theta^{(t+1)} = \arg \max_{\theta} [J(\theta, \theta) + \lambda \text{sup}_{\theta} (\theta)].$$

Therefore

$$\begin{aligned}\text{Lsemi-sup}(\theta^{(t)}) &= \text{Lsup}(\theta^{(t)}) + \lambda \text{sup}_{\theta} (\theta^{(t)}) \\ &= J(\theta^{(t)}, \theta^{(t)}) + \lambda \text{sup}_{\theta} (\theta^{(t)}) \\ &\leq J(\theta^{(t)}, \theta^{(t+1)}) + \lambda \text{sup}_{\theta} (\theta^{(t+1)}) \\ &\leq \text{Lsup}(\theta^{(t+1)}) + \lambda \text{sup}_{\theta} (\theta^{(t+1)}) \\ &= \text{Lsemi-sup}(\theta^{(t+1)})\end{aligned}$$

Hence,  $\text{Lsemi-sup}(\theta^{(t+1)}) \geq \text{Lsemi-sup}(\theta^{(t)})$  and the algorithm will converge monotonically.

d. b.' In the  $t$ -step, we need to re-estimate the latent variables  $Z^{(i)}, g, l = 1, \dots, n$ . We have:

$$\begin{aligned}
 w_j^{(i)} &= P(Z^{(i)} = j | X^{(i)}; \theta) \\
 &= \frac{P(Z^{(i)} = j; \theta) P(X^{(i)} | Z^{(i)} = j; \theta)}{\sum_{l=1}^k P(Z^{(i)} = l; \theta) P(X^{(i)} | Z^{(i)} = l; \theta)} \\
 &= \frac{(t\pi)^{\text{data}} \sum_j \exp\left\{-\frac{1}{2}(X^{(i)} - \mu_j)^T \Sigma_j^{-1} (X^{(i)} - \mu_j)\right\} \phi_j}{\sum_{l=1}^k \frac{1}{(2\pi)^{d/2}} \exp\left\{-\frac{1}{2}(X^{(i)} - \mu_l)^T \Sigma_l^{-1} (X^{(i)} - \mu_l)\right\} \phi_l}.
 \end{aligned}$$

a.c.

In the off-step, we re-estimate the model parameters  $\beta_0$ 's,  $\beta_1$ 's and  $\beta_2$ 's,  $j \in \{1, \dots, k\}$  to maximize the log-likelihood function.

$$\sum_{l=1}^L \sum_{j=1}^{n_l} w_j^{(l)} \log p(x^{(l)}, z^{(l)} = j; \theta) + d \sum_{l=1}^L \log p(x^{(l)}, z^{(l)}; \theta)$$

which is the same as maximizing

$$\sum_{l=1}^L \sum_{j=1}^{n_l} w_j^{(l)} \log p(x^{(l)}, z^{(l)} = j; \theta) + d \sum_{l=1}^L \sum_{j=1}^{n_l} w_j^{(l)} \log p(x^{(l)}, z^{(l)}; \theta)$$

If we append the labeled dataset to the, we get the entire training set of  $n+h$  examples  $(x^{(i)}, z^{(i)})$ , of which the first  $n$  are labeled and the last  $h$  index by

$x^{(1)}, \dots, x^{(n)}$  are labeled. Also about for labeled examples  
 $w_j^{(i)} = b + \sum_{j=1}^k z^{(i)}_j$ ,  $i \in \{1, \dots, n\}$ . So the objective can be  
written as

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log p(x^{(i)}, Z^{(i)} = j; \theta) + \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log p(x^{(i)}, Z^{(i)} = j; \theta) \\ &= \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log p(x^{(i)}, Z^{(i)} = j; \theta) \end{aligned}$$

This is the same as the objective we used to maximize in the 1st-step of the classical Bayes model (excluding additive constants) is equivalent to the one mentioned. By following the approach presented in the lecture notes, we can derive the update rule in the following manner.

$$\phi_j = \frac{1}{n+h} \sum_{i=1}^{n+h} w_j^{(i)}$$

$$\mu_j = \frac{\sum_{i=1}^{n+h} w_j^{(i)} x^{(i)}}{\sum_{i=1}^{n+h} w_j^{(i)}}$$

$$\sum_j = \frac{\sum_{i=1}^{n+h} w_j^{(i)} (x^{(i)} - \mu_j)^T (x^{(i)} - \mu_j)}{\sum_{i=1}^{n+h} w_j^{(i)}}$$

Note that  $w_j^{(i)}$  is computed as conditional probability in the  $t$ -step for  $i \in \{1, \dots, n+h\}$  and is set to  $\alpha \cdot \beta z^{(i)}$  for  $i \in \{n+1, \dots, n+h\}$ .

d. f. i.

i. The unsupervised model takes longer to converge.

"

ii. Due to having no observed data to base its class functions on, it is to be expected that the unsupervised model does not label all data points correctly. The semi-supervised model does not have this problem because it knows the correct labels of at least some samples. But furthermore the unsupervised model can start to create comparable distributions of the classes, while the semi-

"monitored model seems to always end up with the same distribution."

III. The unsupervised model seems to have trouble to identify the high-variance distribution and to distinguish two of the low-variance distributions. The semi-supervised model does seem to find reasonable assignments everywhere.

I.e. In LCF, we want to maximize as a function of  $w$  the following objective:

$$\begin{aligned}
 \ell(w) &= \sum_{i=1}^n \log p_x(x^{(i)}) \\
 &= \sum_{i=1}^n \log \left( \beta_3 \left( w^T x^{(i)} \right) / \|w\| \right) \\
 &= \sum_{i=1}^n \log \left( \frac{1}{(2\pi)^d} \exp \left\{ -\frac{1}{2} (w^T x^{(i)})^T (w^T x^{(i)}) \right\} \right) \\
 &= \sum_{i=1}^n \log \left( \frac{1}{(2\pi)^d} \exp \left\{ -\frac{1}{2} (w^T x^{(i)})^T (w^T x^{(i)}) \right\} \right) \\
 &= \sum_{i=1}^n \left( \frac{1}{2} \log(6\pi) - \frac{1}{2} (x^{(i)})^T w^T w \right) + \log(\|w\|)
 \end{aligned}$$

To maximize this objective, we will compute its gradient and set to 0. We have:

$$\begin{aligned}
 \nabla_W \ell(W) &= \sum_{i=1}^n \left( -\frac{1}{2} \nabla_{W^{(i)}}^T W^{(i)} X^{(i)} + \nabla_W \log |W| \right) \\
 &= \sum_{i=1}^n \left( -W^{(i)} X^{(i)} X^{(i)T} + (W^{-1})^T \right) \\
 &= -W \left( \sum_{i=1}^n X^{(i)} X^{(i)T} \right) + n (W^{-1})^T \\
 &= -W X^T X + n (W^{-1})^T
 \end{aligned}$$

Set this equal to 0, we can get:  $W X^T X = \left(\frac{n}{n} X^T X\right)^{-1}$  assuming that the right-hand side is invertible. Let  $Y = \left(\frac{n}{n} X^T X\right)^{-1}$ ,

then  $\mathbf{Y}$  is positive semi-definite. With  $\mathbf{U}$  and  $\mathbf{V}$  here  
 $\mathbf{W} = \mathbf{U}\Sigma\mathbf{V}^T$  where  $\mathbf{U}$ ,  $\mathbf{V}$  are orthogonal and  $\Sigma$  is diagonal.  
Then, we have

$$\mathbf{W}^T\mathbf{W} = (\mathbf{V}\Sigma\mathbf{U}^T)(\mathbf{U}\Sigma\mathbf{V}^T) = \mathbf{V}\Sigma(\mathbf{U}\mathbf{U}^T)\Sigma\mathbf{V}^T = \mathbf{V}\Sigma^2\mathbf{V}^T = \mathbf{Y}$$

So we can complete the eigendecomposition of  $\mathbf{Y}$  to recover  
 $\Sigma^*$ ,  $\mathbf{V}^*$  and means that  $\mathbf{W} = \mathbf{U}\Sigma^*\mathbf{V}^T$  with an arbitrary  
orthogonal matrix  $\mathbf{U}$ . This arbitrary rotational component  
 $\mathbf{U}$  can be determined from the data  $\mathbf{Y}$  which leads to  
ambiguity. The ICA can not recover the original sources.

3.6.

In general we have

$$(\nabla_w \log g^i(w^T x^{(i)})) V = \frac{g'(w_j^T x^{(i)})}{g''(w_j^T x^{(i)})} e_j^T V x^{(i)}$$

$$= \frac{g'(w_j^T x^{(i)})}{g''(w_j^T x^{(i)})} \text{tr}((x^{(i)})^T e_j^T V)$$

and therefore

$$\begin{aligned} \nabla_w \log g^i(w^T x^{(i)}) &= \frac{g'(w_j^T x^{(i)})}{g''(w_j^T x^{(i)})} \left( (x^{(i)})^T e_j^T \right)^T \\ &= \frac{g'(w_j^T x^{(i)})}{g''(w_j^T x^{(i)})} g^i(x^{(i)})^T \end{aligned}$$

To fit our weights with the gradient of the log-likelihood (for a single example), i.e.  $M = 1$  we:

$$\nabla_w \ell(w) = (w^T)^{-1} + \sum_{j=1}^n \left[ \frac{g'(w_j^T x^{(i)})}{g''(w_j^T x^{(i)})} e_j (x^{(i)})^T \right]$$

$$= (w^T)^{-1} + \begin{bmatrix} \frac{g'(w_1^T x^{(i)})}{g''(w_1^T x^{(i)})} \\ \vdots \\ \frac{g'(w_n^T x^{(i)})}{g''(w_n^T x^{(i)})} \end{bmatrix} (x^{(i)})^T$$

In particular for a Logistic unitarity  $f'(s) = f_1(s) = \frac{1}{2} \exp(-Bs)$   
we get the gradient update rule

$$\nabla_w L(w) = w + \alpha \left[ \begin{matrix} \text{sgn}(w_1^T x^{(i)}) \\ \vdots \\ \text{sgn}(w_n^T x^{(i)}) \end{matrix} \right] \left[ \begin{matrix} (x^{(i)})^T \\ \vdots \\ (x^{(i)})^T \end{matrix} \right].$$

4.a.

Let  $X \in \mathbb{R}^{n \times d}$ . If  $X^T$  is invertible, then  $\beta$  achieves zero cost in Eq (1) iff:

$$\beta = X^T(XX^T)^{-1}y + \zeta$$

for some  $\zeta$  in the subspace orthogonal to all the data (that is, for some  $\zeta$  such that  $Z^T \zeta^{(i)} = 0$ ,  $\forall 1 \leq i \leq n$ ).

Let  $\beta$  be a global minimizer of  $J_F(\beta)$ . Then we have

$$J(\beta) = 0$$

This means that

$$X\beta - y = 0 \quad \text{or} \quad \beta = X^T y$$

New, let  $\zeta$  be any vector such that  $Z^T \zeta^{(i)} = 0 \quad \forall 1 \leq i \leq n$ .

Then we have

$$(B + \gamma I)^T e^{(i)} = \beta^T x^{(i)} + \gamma x^{(i)} = 0$$

for all  $1 \leq i \leq n$ . This means that  $B + \gamma I$  is orthogonal to all the data.

Let  $\beta$  be any vector of the form:

$$\beta = H^T (H H^T)^{-1} y + \gamma$$

for some  $\gamma$  in the subspace orthogonal to all the data.

Then we have

$$J(\beta) = \frac{1}{2} \|H\beta - y\|^2$$

Since  $\gamma$  is orthogonal to all the data, we have

$$\|H\beta - y\|^2 = \|H^T (H H^T)^{-1} y\|^2$$

Now, we know that  $H^T$  is invertible, so we can write

$$H^T(HH^T)^{-1}y = H^T(HH^T)^{-1}(HH^T)y = y$$

This means that

$$\|H\beta - y\|_2 = \|y\|_2 = 0$$

Therefore  $\beta$  is a global minimizer of  $L_F(\beta)$ . And we have shown that  $\beta$  achieves zero loss in  $L_F(\beta)$  iff  $\beta = H^T(HH^T)^{-1}y + \zeta$  for some  $\zeta$  in the subspace orthogonal to all the data. This implies that there is infinite number of  $\beta$ 's such that  $L_F(\beta)$  is minimized.

4. b.

Let  $X \in \mathbb{R}^{n \times d}$  be a matrix with  $n < d$ . If  $X^T$  is invertible, then for any  $\beta$  such that  $J(\beta) = 0$ , we have

$$\|\beta\|_2 \leq \|\beta\|_2 \\ \text{where } \beta = X^T(XX^T)^{-1}y$$

Let  $\beta$  be any vector of the form

$$\beta = X^T(XX^T)^{-1}y + \zeta$$

for some  $\zeta$  in the subspace orthogonal to all data.

then we have

$$\|\beta\|_2^2 = \|\beta\|_2^2 + \|\zeta\|_2^2$$

To see this we have

$$\|\beta\|_2^2 = (\beta)^T \beta = (X^T(XX^T)^{-1}Y + Z)^T (X^T(XX^T)^{-1}Y + Z)$$

Expanding we get

$$\|\beta\|_2^2 = X^T(XX^T)^{-1}Y (X(X^T)^{-1}X\beta + Z^T Z)$$

Since  $Z$  is orthogonal to all the data, we have  $Z^T Z^{(i)} = 0$  for  $1 \leq i \leq n$ . This means that

$$X^T(XX^T)^{-1}Y (X(X^T)^{-1}X\beta + Z^T Z) = X^T(XX^T)^{-1}Y Y = \|Y\|_2^2$$

Therefore we have

$$\|\beta\|_2^2 = \|\beta\|_2^2 + \|Z\|_2^2$$

Let  $\beta$  be any vector such that  $f(\beta) = 0$ . Then we get

$$\|\beta\|_2^2 = \|\rho\|_2^2 + \|\zeta\|_2^2$$

for some  $\zeta$  is the subspace orthogonal to all the data  
line.  $\|\zeta\|_2^2 \geq 0$ , we have

$$\|\beta\|_2^2 \geq \|\rho\|_2^2$$

This means that  $\rho$  is minimum norm solution. And we  
have shown that for any  $\beta$  such that  $f(\beta) = 0$ , we have

$$\|\rho\|_2 \leq \|\beta\|_2$$

4. d

Let  $X \in \mathbb{R}^{n \times d}$  be a matrix with  $n < d$ . If  $X^T$  is invertible, then the gradient descent algorithm with zero initialization always converges to the minimum norm solution.

Let  $B(0) = 0$ . Then we have

$$B(t) = -\frac{1}{h} X^T (X B(0) - y) = -\frac{1}{h} X^T X y$$

Since  $X^T$  invertible, we can write

$$B(t) = -\frac{1}{h} X^T X y = -\frac{1}{h} X^T (X X^T)^{-1} (X^T y) = -\frac{1}{h} X^T y$$

This shows that  $B(t)$  is a linear combination of

$$\begin{pmatrix} 1 \\ X^{(1)} \\ X^{(2)} \\ \vdots \\ X^{(n)} \end{pmatrix}^T y$$

Assume that  $B^{(t)}$  is a linear combination of  $\{x^{(1)}, \dots, x^{(n)}\}$  for

some  $T \geq 1$ . Then we have

$$\begin{aligned}\beta^{(t+1)} &= \beta^{(t)} - \frac{\gamma}{h} X^T (y_{\beta^{(t)}} - f(X^T \beta^{(t)})) \\ &= \beta^{(t)} - \frac{\gamma}{h} X^T (f(\beta^{(t)}) - V^{(t)}).\end{aligned}$$

Since  $\beta^{(t)}$  is a linear combination of  $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ , so is  $\beta^{(t+1)}$ . This means that  $\beta^{(t+1)}$  is also a linear combination of  $f(x^{(1)}), f(x^{(2)}), \dots, f(x^{(n)})$ . By the principle of mathematical induction, we have shown that  $\beta^{(t)}$  is a linear combination of  $x^{(1)}, x^{(2)}, \dots, x^{(n)}$  for all  $t \geq 0$ .

Now let  $\beta^*$  be any solution to the gradient descent<sup>†</sup> with zero initialization that converges to  $J(\beta^*) = 0$

Then we have

$$\beta = \beta^{(t)} \text{ for some } t \geq 0,$$

Since  $\beta^{(t)}$  is a linear combination of  $\beta_1^{(t)}, \beta_2^{(t)}, \dots, \beta_m^{(t)}$ , so  
is  $\beta$ . This means that  $\beta$  is in the subspace orthogonal  
to all the data. But we know that the minimum norm  
solution is also in the subspace orthogonal to all the data.  
Therefore,  $\beta$  must be the minimum norm solution.

We have shown that the gradient descent algorithm  
with zero initialization always converges to the minimum  
norm solution.

4.C.

To prove that there are infinitely many solutions with zero cost in the present quadratic parameterized model, let's first establish that  $y^{(i)}$  is linear in  $\beta^{(i)}$ . We can use the parameter  $\beta^*$  from sub-section (a) to represent  $y^{(i)}$  as follows

$$y^{(i)} = (x^{(i)})^T ((\beta^*)_{01} - (\beta^*)_{02}) \\ = (x^{(i)})^T (\beta^*)^T \\ = (\beta^*)^T x^{(i)}$$

This shows that  $y^{(i)}$  is linear in  $\beta^{(i)}$  with parameter  $\beta^*$ .

Now, let's consider the cost function  $J(\beta)$  for linear model:

$$J(\beta) = \frac{1}{n} \| X\beta - y \|_2^2$$

We can rewrite the quadratically parameterized model  $f_{\theta, \phi}(x)$  in terms of  $\beta$  as follows:

$$f_{\theta, \phi}(x) = x^\top (\phi^{02} - \phi^{01})$$

$$= (\beta^*)^\top x$$

Comparing this with the linear model, we can see that  $\phi^{02}$  and  $\phi^{01}$  can be considered as solutions with zero cost, where  $\ell$  represents the different possible solutions.

Therefore, by mapping the parameters  $\phi$  and  $\ell$  to parameters  $\beta$  in the linear model, we show that there are infinitely many solutions with zero cost in the quadratically parametrized models.

- 4.9 The models with the smaller initialization parameters ( $\delta = 0.01$ , and  $\delta = 0.03$ ) achieve lower training errors. This suggests that smaller initialization parameters can help with overfitting.
- Model with the smallest initialization parameter achieves the best validation error. This suggests that smaller initialization parameters can help to improve generalization performance.

4. i. Typically SGD can find solutions that generalize better than CG due to the noise introduced by the mini-batches. This stochasticity allows the optimizer to explore different areas of the parameter space and avoids getting stuck in poor local optima. But it's not case on finds better solution than SGD because our data set is small. And CG with the smallest initialization parameter generalizes better.