

Recap from previous week

- Can fit different kinds of models to data
 - ↳ PDEs, ODEs, regressors
- Review new HW
 - focus on computation
- Discuss in class exam 3/28
- This week - shifting toward dimension reduction + hidden data

Given inputs $x \in \mathbb{R}^D$, seek mapping

$$x \rightarrow z \in \mathbb{R}^L \quad L \ll D$$

- Typically used for efficiency, generative modeling

ex for images $D = \# \text{ pixels}$

- For physics, we want to reduce DOFs to make simulation orders of magnitude faster
 - Projection-based reduced order models (**PCA**)
 - Koopman operator embeddings (**autoencoder + NODE**)
 - variational inference for generative models

(Office hrs tomorrow)

Principal component analysis (PCA)

Idea: Postulate a linear mapping to

$$\text{encode: } z = w^T x \quad x \rightarrow z$$

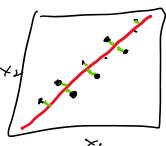
and decode: $\hat{x} = wz$

seek w which minimizes a measure of reconstruction loss

$$f(w) = \frac{1}{N} \sum_{i=1}^N \|x_i - \text{decode}_w(z_i)\|^2$$

Before the math

goal: identifying low dimensional linear subspace

Derivation

Given dataset $\mathcal{D} = \{x_n\}_{n=1}^N$ scaled to have zero mean $\Rightarrow \sum_n x_n = 0$

Assume data can be decomposed into linear combination of basis functions $w_k \in \mathbb{R}^D$, $k=1, \dots, L$

$$x_n \approx \sum_{k=1}^L z_{nk} w_k$$

↑ latent "code"

Note: It will turn out that w_k are the K -largest eigen vectors of covariance matrix Σ

$$\Sigma = Q \Lambda Q^{-1}$$

↑
first
 K columns

$$f(w, z) = \frac{1}{N} \|x - zw^T\|_F^2$$

$$= \frac{1}{N} \sum_n \|x_n - zw_n\|^2$$

Recall Orthogonal Matrix
 $Q^T Q = Q Q^T = I$

Why orthonormal? $\|x - \text{decade } 0 \text{ encode } x\|^2$
 $\geq \|x - w^T w x\|^2$ $\forall x \in \text{span}(w)$
 ≥ 0

Strategy - proceed by mathematical induction

Recall If you want to prove $\{x_n\}_{n=1}^N$ statements are true

- Principle of Mathematical Induction
- ① Prove x_1 is true base case
 - ② Prove x_n is true implies x_{n+1} true

will get a greedy algorithm

Base case

- Find single basis vector $w_i \in \mathbb{R}^D$
Assuming orthonormality $w_i^T w_i = 1$

Let $\tilde{z}_i = [z_{i1}, \dots, z_{iN}]$

$$\begin{aligned} f(w_i, \tilde{z}_i) &= \frac{1}{N} \sum_{n=1}^N \|x_n - z_{ni} w_i\|^2 \\ &= \frac{1}{N} \sum_{n=1}^N (x_n - z_{ni} w_i)^T (x_n - z_{ni} w_i) \\ &= \frac{1}{N} \sum_{n=1}^N x_n^T x_n - 2 z_{ni} w_i^T x_n + z_{ni}^2 w_i^T w_i \end{aligned}$$

At optimal \tilde{z}_i :

$$0 = \frac{\partial}{\partial z_{ni}} f = \frac{1}{N} [-2 w_i^T x_n + 2 z_{ni}]$$

$$\Rightarrow \tilde{z}_{ni}^* = w_i^T x_n$$

Plugging back into loss

$$\begin{aligned} f(w_i, \tilde{z}_i^*) &= \frac{1}{N} \sum_n x_n^T x_n - 2 \tilde{z}_{ni}^2 + \tilde{z}_{ni}^2 \\ &= \frac{1}{N} \sum_n [x_n^T x_n - z_{ni}^2] \\ &\approx \text{const} - \frac{1}{N} \sum_n z_{ni}^2 \\ &= \text{const} - \frac{1}{N} \sum_n w_i^T x_n x_n^T w_i \end{aligned}$$

$$f(w_i) = -w_i^T \sum_n^N w_i$$

Recall that $\hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2$ but we scaled data

To find optimal unit normal basis

$$\max_w w_i^T \sum_n^N w_i$$

$$\text{s.t. } \|w_i\| = 1$$

Lagrange Multipliers

$$0 = \partial_{w_i} \left[w_i^T \sum_n^N w_i - \lambda_1 (w_i^T w_i - 1) \right]$$

$$0 = 2 \sum_n^N w_i - \lambda_1 2 w_i$$

$$\Rightarrow \sum_n^N w_i = \lambda_1 w_i \quad \text{eigenvalue equation}$$

$$\therefore w^T \sum_n^N w_i = \lambda_1 w_i^T w_i$$

$$= \lambda_1, \quad \text{takes largest value at biggest eigenvalue}$$

Proof by induction

Given w_1, \dots, w_{i-1} , seek w_i satisfying
 orthonormal conditions $\begin{cases} (1) w_i^T w_j = 0 & \text{for } j < i \\ (2) w_i^T w_i = 1 \end{cases}$

$$\begin{aligned} \mathcal{L}(w_1, \dots, w_{i-1}, w_i, z_1, \dots, z_{i-1}, z_i) &= \\ &= \frac{1}{N} \sum_{n=1}^N \|x_n - \sum_{j=1}^{i-1} z_{nj} w_j - z_{ni} w_i\|^2 \end{aligned}$$

The $\frac{\partial}{\partial w_j} \mathcal{L} = 0$ for $j < i$, we recover the minimizer from before
 $\frac{\partial}{\partial z_{ni}} \mathcal{L} = 0$ $z_{nj} = w_j^T x_n$
 w_i = i th eigenvector

Fixing $j < i$ terms in loss we can show

$$0 = \frac{\partial}{\partial z_i} \mathcal{L} \Rightarrow z_{ni} = w_i^T x_n$$

$$\mathcal{L}(w_i) = \text{const} - w_i^T \sum_{j=1}^i w_j$$

$$\max_{w_i} \mathcal{L}(w_i)$$

$$\text{s.t. } w_i^T w_i = 1$$

$$w_i^T w_j = 0 \quad \text{for all } j < i$$

$$\sum_{j=1}^i w_j = \lambda_i w_i \quad \text{why?}$$

biggest is

largest λ_1

but $w_i^T w_1 = 0$

next is λ_2

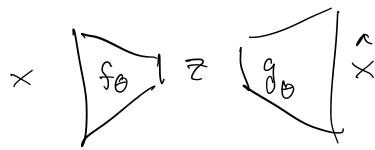
$$w_i^T w_2 \geq 0$$

etc

Guarantees

- orthonormal basis
- complete recovery

Nonlinear encoding



- Trivial to set up
- No guarantees
- No "data" for what a "good" z would be

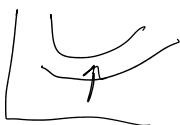
Hidden Data

Murphy § 8.7.2

- Consider $D = \{(y_i, z_i)\}_{i=1}^N$
 - observable
 - hidden
- log likelihood is not possible, because we have no values to plug in for z in $\log p(y, z)$
 - unknown
- Define observed log likelihood
$$LL(\theta) = \sum_{n=1}^N \log(p(y_n, z_n | \theta))$$
- Expectation Maximization gets around this by deriving an alternative loss that bounds $LL(\theta)$ from below

sketch

$$F_\theta < LL_\theta$$



$$\max F_\theta < \max LL_\theta$$

and if the bound is tight they are equivalent

Recall

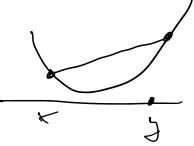
Jensen's inequality

For convex $g(x)$ $E[g(x)] \geq g(E[x])$

ex $g(x) = x^2$

formally
$$g(\theta x + (1-\theta)y) \leq \theta g(x) + (1-\theta)g(y)$$

 $\forall \theta \in [0, 1]$



Derivation

$$LL(\theta) = \sum_{n=1}^N \log \sum_{z_n} p(y_n, z_n | \theta)$$

$$\begin{aligned} \text{for arbitrary dist. } f_n(z) &= \sum_n \log \sum_{z_n} g_n(z_n) \frac{p(y_n, z_n | \theta)}{g_n(z_n)} \\ \text{think of as secret dist.} &= \sum_n \log \mathbb{E}_{z_n \sim g} \left[\frac{p(y_n, z_n | \theta)}{g_n(z_n)} \right] \\ &\geq \sum_n \mathbb{E}_{z_n \sim g} \left[\log \frac{p(y_n, z_n | \theta)}{g_n(z_n)} \right] \\ &= \sum_n \sum_{z_n} g_n(z_n) \log \frac{p(y_n, z_n | \theta)}{g_n(z_n)} \\ &= \sum_n \mathbb{E}_{z_n \sim g_n} \left[\log p(y_n, z_n | \theta) \right] + H(g_n) \\ &= \sum_n \underbrace{\mathcal{E}(\theta, g_n | y_n)}_{\text{Evidence Lower Bound (ELBO)}} \end{aligned}$$

entropy
 $H(f) = \sum_n f_n \log f_n$

We can express the ELBO

$$\begin{aligned} \mathcal{E}(\theta, g_n | y_n) &= \sum_{z_n} g_n(z_n) \log \frac{p(y_n, z_n | \theta)}{g_n(z_n)} \\ &= \sum_{z_n} g_n(z_n) \log \underbrace{\frac{p(z_n | y_n, \theta) p(y_n | \theta)}{g_n(z_n)}}_{p(z_n | y_n, \theta)} \end{aligned}$$

$$\begin{aligned} \mathcal{E}(\theta, g_n | y_n) &= \sum_{z_n} g_n \log \frac{p(z_n | y_n, \theta)}{g_n} - \sum_{z_n} g_n(z_n) \log p(y_n | \theta) \\ &= -KL(g_n(z_n) || p(z_n | y_n, \theta)) + \log p(y_n | \theta) \end{aligned}$$

① ②

Recall KL divergence

$$\begin{aligned} KL(f || g) &> 0 \\ KL(f || g) = 0 &\Rightarrow f = g \end{aligned}$$

maximizing E wrt g , set $g_n^* = p(z_n | y_n, \theta)$

- (why?)
- (2) indep of z_n
- (1) is least negative when $= 0$

Plug back in

$$E(\theta, g^* | y_n) = \log p(y_n | \theta) \quad \text{i.e. KL term drops out}$$

And maximizing ELBO w/ choice

$$g_n(z_n) = p(z_n | y_n, \theta) \quad \begin{matrix} \uparrow \\ \text{Posterior} \end{matrix} \quad \text{is tight!} \quad \begin{matrix} \uparrow \\ \text{Dist.} \end{matrix}$$

end of lecture 2/20

EM - Algorithm

Given initial θ^0 params

for $t = 1, \dots, N$ iterations :

$$g_n^t = p(z_n | y_n, \theta^{t-1}) \quad - \underline{\text{E-step}}$$

$$\theta^t = \arg \max_{\theta} E_{z_n \sim g_n^t} [\log p(y_n, z_n | \theta)] \quad - \underline{\text{M-step}}$$

Example - Gaussian Mixture

Have data y_n but no labels for

Have data y_n but no labels for cluster ownership

$$p(y) = \sum_{c=1}^{N_{\text{cluster}}} p(y|z) p(z=c) \underbrace{\pi_{\text{cat}}(t)}_{N(y; \mu_c, \Sigma_c)}$$

To derive E-step

$\hat{p}_{y|z}$

$$\begin{aligned} p(z=c | y_d) &= \frac{p(y_d | z=c) p(z=c)}{p(y_d)} \\ &= \frac{p(y_d | z=c) p(z=c)}{\sum_{c'} p(y_d | z=c') p(z=c')} \end{aligned}$$

law of total probability

Here, interpret
 $\theta = \{\mu_c, \Sigma_c\}_{c=1}^{N_{\text{cluster}}}$

$$q_d^t(z_d) = \frac{\pi_c N(y_d; \mu_c, \Sigma_c)}{\sum_{c'} \pi_{c'} N(y_d; \mu_{c'}, \Sigma_{c'})}$$

M-step

$$f^t = \sum_d \mathbb{E}_{q_d^t(z_d)} \left[\log p(y_d, z_d | \theta) \right]$$

$$= \sum_d \sum_c q_d^t(z_d=c) \left(\log p(y_d | z_d=c, \theta) p(z_d=c | \theta) \right)$$

$$= \sum_s \sum_c q_d^t(z_d=c) \left[\log N(y_d | \mu_c, \Sigma_c) + \log(\pi_c) \right]$$

$$= \sum_d \sum_c g_d^t(z_d=c) \left[\log N(y_d | \mu_c, \Sigma_c) + \log(\pi_c) \right]$$

Nice things

$$\frac{\partial L^t}{\partial \pi_c} = 0 \Rightarrow \pi_c^{t+1} = \sum_d \frac{g_d^t(z_d=c)}{N}$$

$$\frac{\partial L^t}{\partial \mu_c} = 0 \Rightarrow \mu_c^{t+1} = \frac{\sum_d g_d^t(z_d=c) y_d}{\sum_d g_d^t(z_d=c)}$$

$$\frac{\partial L^t}{\partial \Sigma_c} = 0 \Rightarrow \Sigma_c^{t+1} = \frac{\sum g_d^t(z_d=c) y_d y_d^T}{\sum g_d^t(z_d=c)}$$

$$- (\mu^{t+1})^T \mu^{t+1}$$

Exercise (future HW)

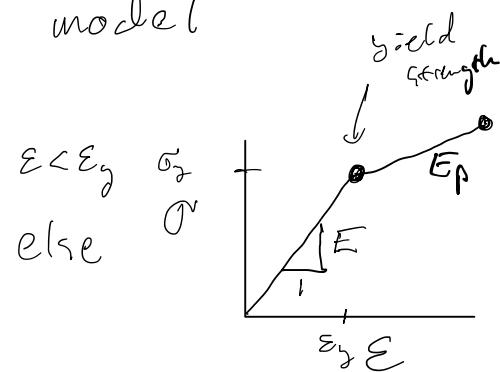
Develop a mixture of experts model which partitions a population curves into N distinct populations. Train w/ EM and use $p(z|x)$ to estimate probability a given curve belongs to a given cluster.

• Step 1 - overall model

Step 1

- Develop expert model
For linear strain hardening model

$$\sigma(\varepsilon | z=c) = \begin{cases} E^c \varepsilon & \varepsilon < \varepsilon_y \\ \sigma_y^c + E_p^c (\varepsilon - \varepsilon_y^c) & \text{else} \end{cases}$$



- Define $p(y | z=c, \varepsilon) = N(y; M=\sigma(\varepsilon | z=c), \sigma^2_c)$

- Define $p(y, z | \varepsilon) = p(y | z=c) p(z=c) \underset{\text{Cat}(\pi)}{\sim}$

E-step $g_d(z_d=c) = \frac{p(y_d | z_d=c) p(z_d=c)}{\sum_{c'} p(y_d | z_d=c') p(z_d=c')}$

M-step $\mathcal{L} = \sum_d \sum_c g_d(z_d=c) \left[\log p(y_d | z_d=c, \varepsilon) + \log p(z_d=c) \right]$

Note : During M-step, because we hold g const
the derivative wrt model params $\underbrace{E^c, E_p^c, \sigma_y^c}_{\theta}$
Gives

$$\theta = \frac{\partial \mathcal{L}}{\partial \theta_c} = \sum_{\theta_c} \sum_c g_d(z_d=c) \log p(y_d | z_d=c, \varepsilon)$$

$$= \sum_c \left[\Gamma_1 + \Gamma_2 + \sum_d (y_d - M_d^c)^2 \right]$$

$$= \mathcal{D}_G' \sum_d^C \sum_c \left[\frac{1}{2} \log \sigma_c^2 + \frac{\sum_e (y_d - M_d^c)^2}{\sigma_c^2} \right]$$

$$\mathcal{O} = \mathcal{D}_G' \sum_d \sum_e q_d^c (y_d - M_d^c)^2$$

★ Decoupled Least squares!

Variational Auto encoders

A related idea to EM, for dim reduction

Recall

Monte Carlo sampler

"Autoencoding Variational Bayes"

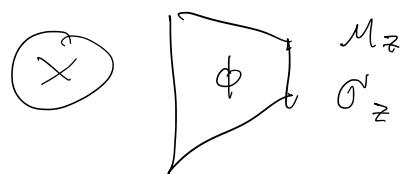
Kingma + Welling
§ 20.3-5 of Murphy

$$\mathbb{E}_{x \sim g} [f(x)] = \frac{1}{N} \sum_i f(x_i) \quad \text{for } x_i \sim g$$

$$\lim_{N \rightarrow \infty} \left| \mathbb{E}_{x \sim g} [f(x)] - \frac{1}{N} \sum_i f(x_i) \right| = 0$$

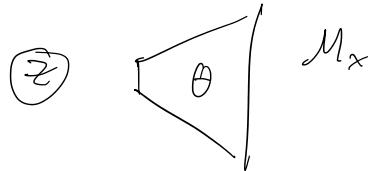
Goal: Replace auto encoder with
~ probabilistic embedding

$$\text{Consider } g(z|x, \phi) = N(\mu_z(x), \sigma_z^2 I)$$



$$p(x|z, \theta) = N(M_x(z), \sigma_x^2)$$

↗ constant noise



ELBO

$$\mathcal{E}(\theta, \phi | x) = \mathbb{E}_{g_\phi(z|x)} \left[\log p_\theta(x|z) - \log g_\phi(z|x) \right]$$

can manipulate

$$= \mathbb{E}_{g_\phi(z|x)} \left[\log p_\theta(x|z) \right] - KL(g_\phi(z|x) || p(z))$$

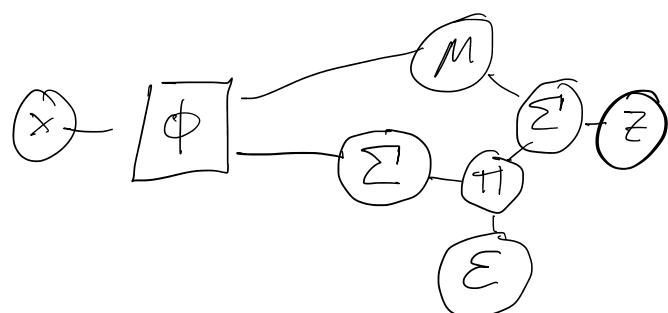
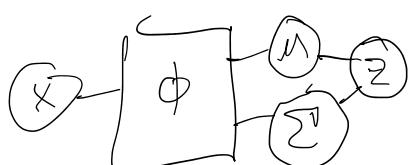
↑
sample w/ monte-carlo

↗ a prior distribution on latent distribution

Reparameterization trick

Sample $\epsilon \sim N(0, I)$

$$\text{Then } z = M_z(x) + \sqrt{\Sigma_z(x)} \epsilon \sim g_\phi(z|x)$$



Mathematically we write this

$$\mathbb{E}_{g(z|x)} [f(z)] = \mathbb{E}_{\varepsilon \sim N(0,1)} [f(m + \sqrt{\varepsilon} z)]$$