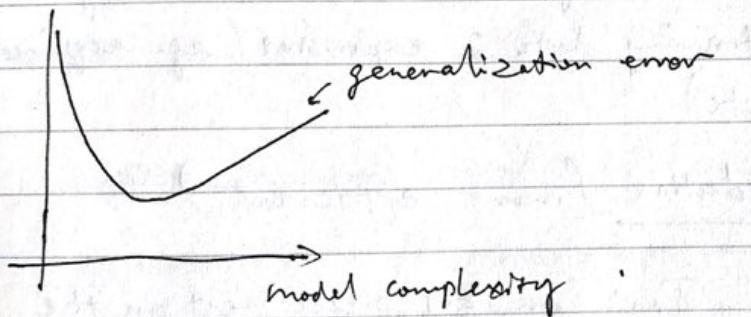


Model Selection



① $\theta_0 + \theta_1 x$

② $\theta_0 + \theta_1 x + \theta_2 x^2$

\vdots

$\theta_0 + \theta_1 x + \dots + \theta_n x^n$

② T - bandwidth parameter in LWR

③ C - in SVM

$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$

$M = \{M_1, M_2, M_3, \dots\}$

$\theta_0 + \theta_1 x$

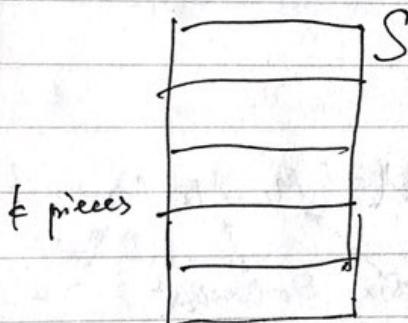
$\theta_0 + \theta_1 x + \theta_2 x^2$

Hold-out cross validation

- Split S into S_{train} (70%)
 S_{cv} (30%) Randomly commonly
- Train each model on S_{train} , test on S_{cv} .
- Pick model with low error on S_{cv} .
- ① Retrain 100% S on best model, or
- ② Output the best model directly.

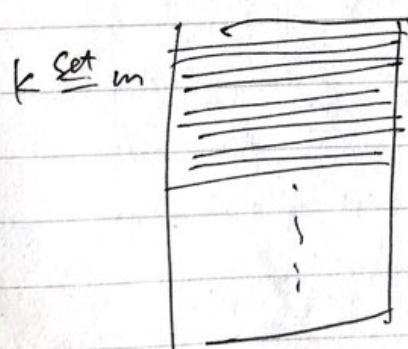
In practice, size of training data is small in many cases.
Because every training data is expensive (e.g. acquired by medical experiments).

k -fold cross validation (more efficient)



- Train on $k-1$ pieces, test on the remaining pieces (do k times totally)
 - Average over the k results
 - Retraining the selected best model
- * $k=10$ is very common

disadvantage: much computationally expensive



Leave-one-out C. V.

(more efficient on using data)

much more computationally expensive
Preferred when only have little training set
(e.g. $m=15$)

People in Structured Risk Minimization (SRM) propose choose k depending on VC-dim, But Andrew Ng personally tend NOT to do that.

→ The upper bounds which proved in learning theory is very loose although in worst case of many models. (for any probability distribution). $(x, y) \in \mathcal{D}$

If plug-in constants into the formula for Learning Theory bounds, you often get extremely optimistic estimates with # examples you need.

Eg.: logistic regression model: if have n or $n+1$ parameters # examples is 10 times # parameters (or tiny times) probably in good shape (also probability).

Feature Selection

High dim. feature space. Eg: 10000 features

(risk of overfitting)

(because of too many ~~parameters~~ parameters in the model)

"relevant" features is much less.

With n features, there are 2^n possible subsets!

Forward search:

- Start with $\mathcal{F} = \emptyset$
- Repeat {
 - (1) For $i=1, \dots, n$: try adding feature i into \mathcal{F} . evaluating model using cross-validation.
 - (2) Set $\mathcal{F} := \mathcal{F} \cup \{ \text{best feature found in (1)} \}$}
- Output best hypothesis found.

→ e.g.: logistic regression model; you already have 100 parameters
you may only need few features.

↓
Terminate the loop when #iter exceeded 100.

"Wrapper" feature selections (computationally expensive)

But tend to work well!

Backward Search:

- Start with $\mathcal{F} = \{1, \dots, n\}$.
- Delete features one at a time.

e.g.: 100 training example.

10000 features. Backward Search is meaningless.

It's NP-Hard to get truly best features.

All above is heuristic search.

52 * For text classification, Do Forward Search is a heavily suffer? too many parameters/features

"Filter" method

For each feature i : Compute some measure of how informative x_i is about y .

(e.g. correlation between x_i & y)

" $\text{Corr}(x_i, y)$ "

$$\text{MI}(x_i, y) = \sum_{\substack{x_i \in \{0,1\} \\ \text{Mutual}}} \sum_{y \in \{0,1\}} p(x_i, y) \log \frac{p(x_i, y)}{p(x_i) p(y)}$$

estimate from training data

[In problem set].

$$= KL(p(x_i, y) || p(x_i) p(y))$$

"Measure of dependence"

Pick the top k features (and choose k using C.V.)

Lecture 11

Bayesian Statistics & Regularization

Degression: Online Learning

Advice for apply ML algorithms.

Prevent overfit via regularization and keep all the ~~parameters~~ ^{parameters}

Linear Regression:

maximum likelihood "Frequentist"

$$\max_{\theta} \prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta)$$

"Bayesian"

$p(\theta)$ - prior e.g. $\theta \sim N(0, \tau^2 I)$

$$S = \{(x^{(i)}, y^{(i)})\}_{i=1}^m$$

calculate:
 $P(\theta|S) \propto \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \right) p(\theta)$

↑
posterior ↑

To make a new prediction on x :

$$P(y|x, S) = \int p(y|x, \theta) p(\theta|S) d\theta \rightarrow$$

Treating θ
as random
variable.

$$\mathbb{E}[y|x, S] = \int_y y p(y|x, S) dy$$

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta|S)$$

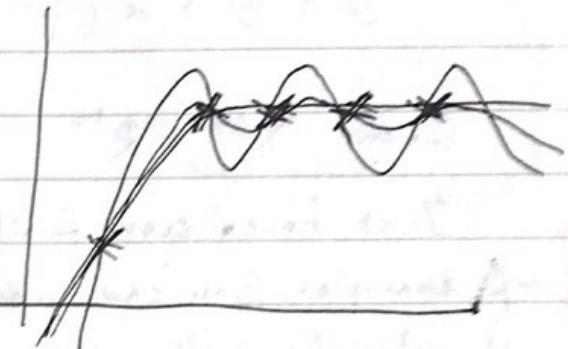
$$= \arg \max_{\theta} \left(\prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \right) p(\theta).$$

To make a prediction:

$$\text{Predict } h_{\hat{\theta}_{MAP}}(x) = \hat{\theta}_{MAP}^T x.$$

$\theta \sim N(0, \tau^2 I)$ ← Set $\theta_0 = 0$: same as delimitate 5th feature
 ↳ cause the effect of smaller parameter & smoother curve.

- ↳ τ smaller and smaller (closer to 0)
 - ↳ curves smoother and smoother
 - ↳ overfit less
- [problem set 3]



e.g.: LR.

$$\min_{\theta} \sum \|y^{(i)} - \theta^T x^{(i)}\|^2$$

$$+ \lambda \|\theta\|^2$$

→ make choose smaller parameter.

* With Bayesian regularization with Gaussian prior,
 logistic regression become very effective.

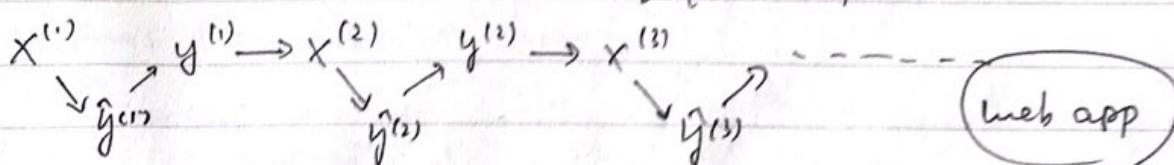
* text classification algorithm..

You can use cross-validation to choose τ^2 . ?

Online Learning

(batch learning algo. so far)

different from



Total Online error:

$$\sum_{i=1}^m \mathbb{I}\{\hat{y}^{(i)} \neq y^{(i)}\}$$

Stochastic grad. desc. algo.

E.g.: Perceptron algorithm

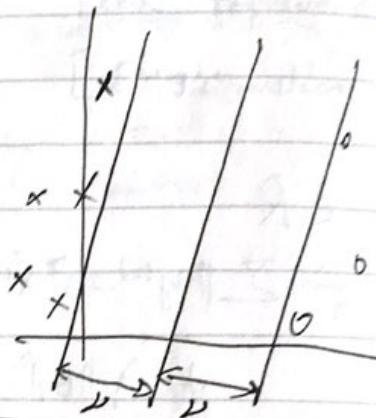
- Initialization: $\Theta := 0$

- After i^{th} example, update

$$\Theta := \Theta + \alpha (y^{(i)} - h_{\Theta}(x^{(i)})) x^{(i)}$$

Given $x^{(i)} \in \mathbb{R}^n$

Just have seen finite number
of examples, you can prove that
the classifier will converge to
perfect separator.



Lecture notes online gives proof (One page)

"Advice for applying Machine Learning" [PDF]

experience

Lecture 12

Unsupervised Learning

- Clustering (k-means)
- Mixture of Gaussians
- Jensen's inequality
- EM (Expectation-Maximization)



Application in biology: cluster genes

"discover structure in the data"

market: cluster customers

news.google.com: cluster news

image segmentation: cluster pixels

K-means algorithm

Input: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \quad x^{(i)} \in \mathbb{R}^n$

1. Initialize cluster centroids

$$\mu_1, \dots, \mu_k \in \mathbb{R}^n$$

2. Repeat until convergence:

i) Set $c^{(i)} = \operatorname{argmin}_j \|x^{(i)} - \mu_j\|$

$$\text{ii) } \mu_j := \frac{\sum_{i=1}^m \mathbb{1}\{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{c^{(i)} = j\}}$$

① K-means is guaranteed to converge

Distortion function

$$J(c, \mu) = \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

K-means is coordinate ascent/descent on J .

② How to choose # clusters ($= k$)?

Randomly choose, pick the best choose.

2 or 4 clusters?
obscure...



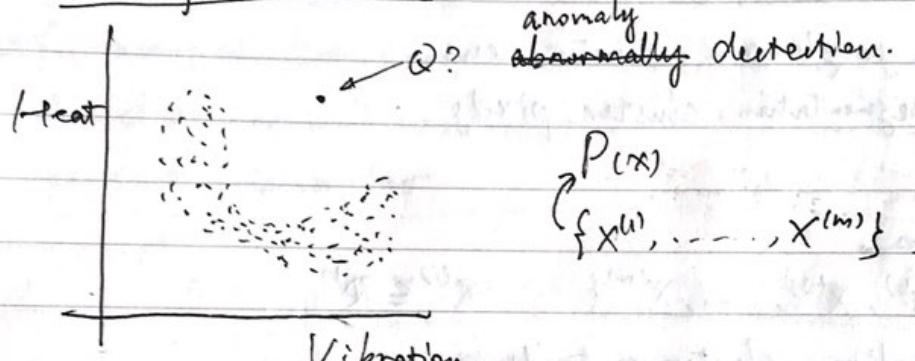
$J(c, \mu)$ is not convex.

③ So k-means may converge to local opt. $\min.$

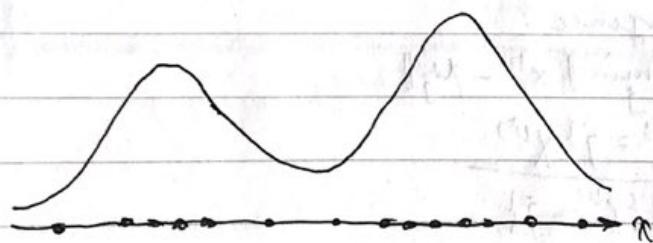
try a few times, choose the least case.

~~1 bus or 2 steps~~
4 stamps or 2 footprints

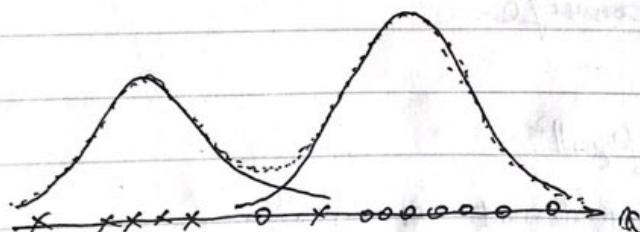
Density estimation



(Aircraft Engine)



We can get the density function without make clear which two Gaussian separately.



There's a latent (hidden/unobserved) random variable z ,

And $x^{(i)}, z^{(i)}$ have a joint distribution

$$P(x^{(i)}, z^{(i)}) = P(x^{(i)}|z^{(i)}) P(z^{(i)})$$

$z^{(i)} \sim \text{Multinomial}(\phi)$ (Bernoulli for 2 Gaussians)
 $(\phi_j \geq 0, \sum_j \phi_j = 1)$

$$x^{(i)}|z^{(i)}=j \sim N(\mu_j, \Sigma_j)$$

Replace $y^{(i)}$ in supervised learning with latent variable $z^{(i)}$
in unsupervised learning.

If we knew $z^{(i)}$'s.

$$\ell(\phi, \mu, \Sigma) = \sum_{i=1}^m \log p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)$$

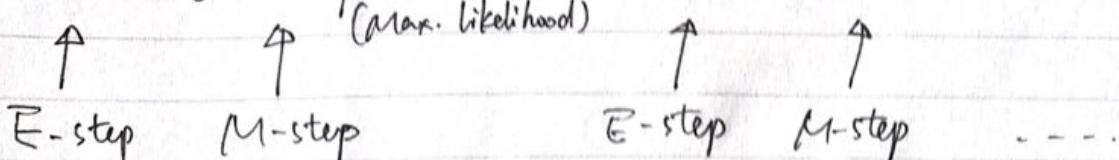
$$\sum_{i=1}^m \log p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma).$$

$$\left\{ \begin{array}{l} \phi_j = \frac{1}{m} \sum_{i=1}^m I\{z^{(i)}=j\} \\ \mu_j = \frac{\sum_{i=1}^m I\{z^{(i)}=j\} x^{(i)}}{\sum_{i=1}^m I\{z^{(i)}=j\}} \end{array} \right.$$

$$\Sigma_j = \dots \quad \text{etc. ...}$$

Guess $z^{(i)}$: bootstrap procedure.

Guess \rightarrow Estimate parameter \rightarrow Guess \rightarrow Estimate \rightarrow Guess $\rightarrow \dots$



EM

Repeat?

E-step: Guess values of $z^{(i)}$'s

$$\text{Let } w_j^{(i)} := P(z^{(i)}=j | x^{(i)}; \phi, \mu, \Sigma)$$

$$= \frac{P(x^{(i)} | z^{(i)}=j) P(z^{(i)}=j)}{\sum_{l=1}^K P(x^{(i)} | z^{(i)}=l) P(z^{(i)}=l)}$$

$$= \frac{\frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left\{ \frac{1}{2} (x^{(i)} - \mu_j^{(i)})^T \Sigma_j^{-1} (x^{(i)} - \mu_j^{(i)}) \right\} \cdot \phi_j}{\sum_{l=1}^K \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_l|^{\frac{1}{2}}} \exp \left\{ \frac{1}{2} (x^{(i)} - \mu_l^{(i)})^T \Sigma_l^{-1} (x^{(i)} - \mu_l^{(i)}) \right\} \cdot \phi_l}$$

M-step: Estimate parameters of Gaussian ℓ

$$\phi_j := \frac{1}{m} \sum_{i=1}^m w_j^{(i)} \quad \text{(GDA)}$$

Compared to Gaussian Discriminant Analysis.

Hence use the probability $x^{(i)}$ comes from Gaussian j .

instead of the indicator whether $x^{(i)}$ comes from Gaussian j .

$$\mu_j := \frac{\sum_{i=1}^m w_j^{(i)} x^{(i)}}{\sum_{i=1}^m w_j^{(i)}}$$

$$\Sigma_j := \frac{\sum_{i=1}^m w_j^{(i)} (x^{(i)} - \mu_j^{(i)}) (x^{(i)} - \mu_j^{(i)})^T}{\sum_{i=1}^m w_j^{(i)}} \quad (\text{i.e. in GDA, } w_j^{(i)} = "1" \text{ or } 0 = I\{y^{(i)}=j\})$$

Jensen's inequality

Let f be a convex function (e.g. $f''(x) \geq 0$)

Let X be a random variable.

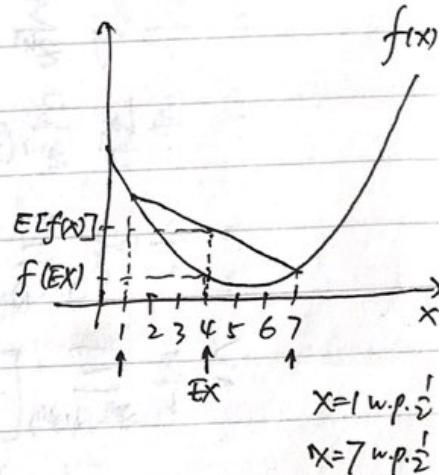
Then $f(\mathbb{E}X) \leq \mathbb{E}[f(X)]$.

Further, if $f''(x) > 0$ (f is strictly convex)
 then $\mathbb{E}[f(X)] = f(\mathbb{E}X)$

$$\Leftrightarrow X = \mathbb{E}X \text{ w.p. } 1$$

If $f'' \leq 0$ (f is concave)

then $f(\mathbb{E}X) \geq \mathbb{E}[f(X)]$ etc.

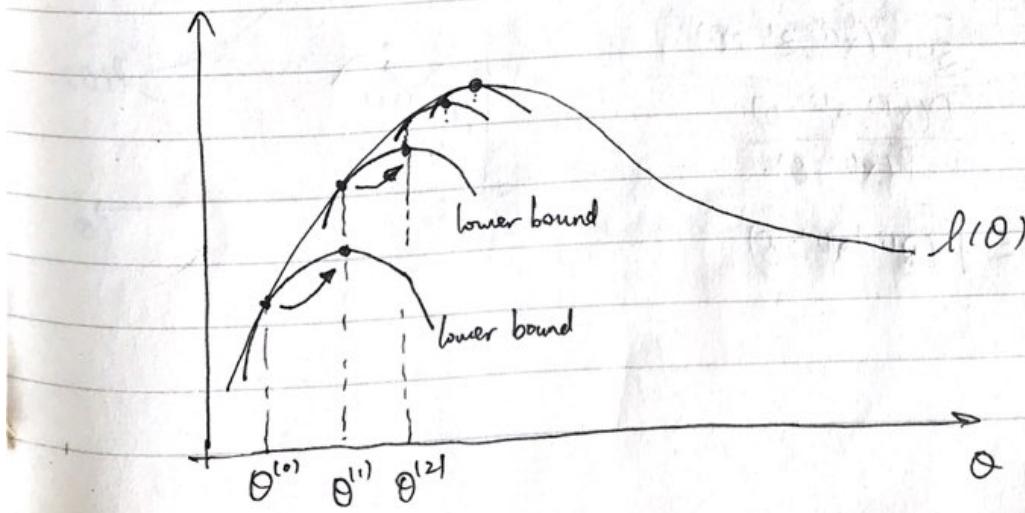


Leave some model for $p(x, z; \theta)$

Observe only X .

Want to maximize

$$\begin{aligned}\ell(\theta) &= \sum_{i=1}^m \log p(x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \left(\sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta) \right) \\ &= \sum_{i=1}^m \log \left(\sum_{l=1}^k p(x^{(i)} | z^{(i)}=l; \theta) p(z^{(i)}=l) \right)\end{aligned}$$



$$\max_{\theta} \sum_i \log P(x^{(i)}; \theta)$$

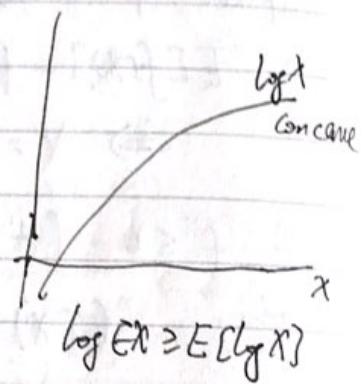
$$= \sum_i \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)$$

$$= \sum_i \log \frac{\sum_{z^{(i)}} Q_i(z^{(i)})}{Q_i(z^{(i)})} \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \quad [Q_i(z^{(i)}) \geq 0, \sum_{z^{(i)}} Q_i(z^{(i)}) = 1]$$

$$= \sum_i \log \mathbb{E}_{z^{(i)} \sim Q(z^{(i)})} \left[\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

$$\geq \sum_i \mathbb{E}_{z^{(i)} \sim Q(z^{(i)})} \left[\log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

$$= \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$



You should make sure it is concave. Many case we met it's concave at all.

Want: $\frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} = \text{constant}$ (for all values of $z^{(i)}$)

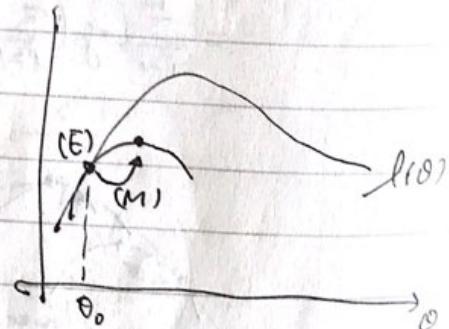
So set: $Q_i(z^{(i)}) \propto P(x^{(i)}, z^{(i)}; \theta)$

$$\sum_{z^{(i)}} Q_i(z^{(i)}) = 1$$

$$Q_i(z^{(i)}) = \frac{P(x^{(i)}, z^{(i)}; \theta)}{\sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)}$$

$$= \frac{P(x^{(i)}, z^{(i)}; \theta)}{P(x^{(i)}; \theta)}$$

$$= P(z^{(i)} | x^{(i)}; \theta)$$



EM algorithm:

E-step:

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

M-step:

$$\theta := \underset{\theta}{\operatorname{argmax}} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

Lecture 13

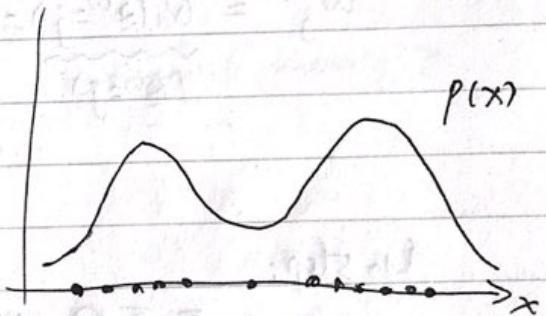
EM

- Mixture of Gaussian
- Mixture of Naive Bayes
- Factor analysis
- Digression: Gaussians

$$\{x^{(1)}, \dots, x^{(m)}\}$$

$$\max_{\theta} \sum_i \log P(x^{(i)}; \theta)$$

$$\max_{\theta} \sum_i \log \sum_{z^{(i)}} P(x^{(i)}, z^{(i)}; \theta)$$



$$P(x) = \sum_z P(x|z) P(z)$$

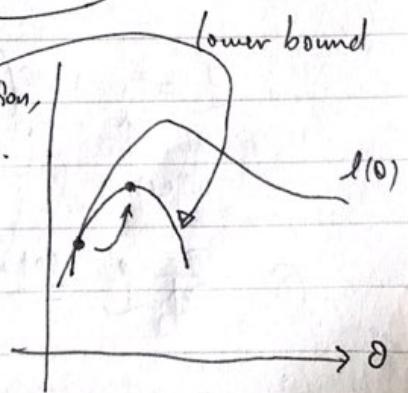
$$z \sim \text{Multinomial}(\phi)$$

$$x|z \sim \mathcal{N}(\mu_j, \Sigma_j)$$

$$\text{E-step: } Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

$$\text{M-step: } \theta := \underset{\theta}{\operatorname{argmax}} \underbrace{\sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}}$$

Can easily get closed-form solution,
rather than $l(\theta)$ which hard to get.



$$\text{Define } J(\theta, Q) = \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

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

Co-ordinate ascent on J.

E-step: Maximize J w.r.t. Q

M-step: Maximize J w.r.t. θ

MoG (Mixture of Gaussian)

E-step:

$$w_j^{(i)} = Q_i(z^{(i)}=j) = P(z^{(i)}=j | x^{(i)}, \phi, \mu, \Sigma)$$

$\stackrel{\text{P}(z^{(i)}=j)}{\longrightarrow}$

$$= \frac{(P(x^{(i)}|z^{(i)}=j) P(z^{(i)}=j))}{\sum_k P(x^{(i)}|z^{(i)}=k) P(z^{(i)}=k)}$$

$x^{(i)} | z^{(i)} \sim N(\mu_j, \Sigma)$
 $z^{(i)} \sim \text{Multinomial}$

M-step:

$$\begin{aligned} & \max_{\phi, \mu, \Sigma} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} \\ &= \sum_i \sum_j w_j^{(i)} \log \frac{\frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j) \right\} \cdot \phi_j}{Q_i(z^{(i)}=j)} \end{aligned}$$

$$\textcircled{1} \quad \nabla_{\mu_j} (\dots) \stackrel{\text{set}}{=} 0 \Rightarrow \mu_j = \frac{\sum w_j^{(i)} x^{(i)}}{\sum w_j^{(i)}}$$

$$\mathcal{L} = (\dots) + \beta \left(\sum_j \phi_j - 1 \right) \quad \left(\sum_j \phi_j = 1 \quad \phi_j \geq 0 \right)$$

↑
objective
Lagrange Multipliers
ignore

$$\textcircled{2} \quad \frac{\partial}{\partial \phi_j} \mathcal{L} \stackrel{\text{set}}{=} 0 \Rightarrow \phi_j = \frac{\sum w_j^{(i)}}{m}$$

$$\textcircled{3} \quad \dots \Rightarrow \Sigma_j = \dots$$

Text clustering: Mixture of Naive Bayes

Training set $\{x^{(1)}, \dots, x^{(m)}\}$. $x^{(i)} \in \{0, 1\}$ $x_j^{(i)} = 1 \{ \text{word } j \text{ appears in doc } i \}$
 $z^{(i)} \in \{0, 1\}$ (2 clusters)

$$z^{(i)} \sim \text{Bernoulli}(\phi)$$

$$P(x^{(i)} | z^{(i)}) = \prod_{j=1}^n P(x_j^{(i)} | z^{(i)})$$

$$P(x_j^{(i)} = 1 | z^{(i)} = 0) = \phi_j | z=0$$

E-step:

$$w^{(i)} = P(z^{(i)} = 1 | x^{(i)}, \phi_j | z, \phi_k) \quad \text{Initialize } \phi \text{ randomly ?!}$$

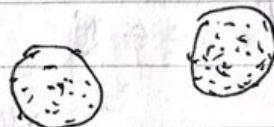
M-step:

$$\phi_j | z^{(i)} = 1 = \frac{\sum_{i=1}^m w^{(i)} I\{x_j^{(i)} = 1\}}{\sum_{i=1}^m w^{(i)}}$$

$$\phi_j | z^{(i)} = 0 = \frac{\sum_{i=1}^m (1 - w^{(i)}) I\{x_j^{(i)} = 1\}}{\sum_{i=1}^m (1 - w^{(i)})}$$

$$\phi_{z^{(i)}} = \frac{\sum_{i=1}^m w^{(i)}}{m}$$

MoG applied: $m \gg n$
 e.g. $n=2$
 $m=100$



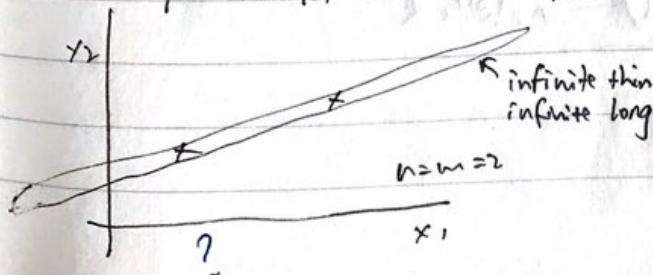
Factor Analysis

What about $n \approx m$, $n \gg m$?

Q $\{x^{(1)}, \dots, x^{(m)}\}$ estimate $P(x)$.

$$x \sim N(\mu, \Sigma) \quad \Sigma \in \mathbb{R}^{n \times n}$$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T \quad \xrightarrow{n \gg m} \text{singular}$$



$$\frac{1}{(2\pi)^{\frac{m}{2}/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$|\Sigma| = 0 !!$$

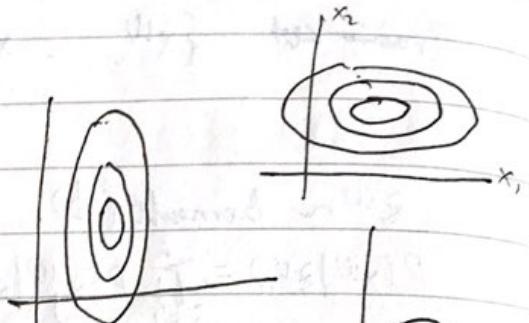
$$P(x) = ?$$

Σ^{-1} undefined!
 Not good model!

② conform Σ to be diagonal $x \sim N(\mu, \Sigma)$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{pmatrix} \in \mathbb{R}^{n \times n}$$

$$\sigma_j^2 = \frac{1}{m} \sum_i (x_j^{(i)} - \mu_j)^2$$



Throw the correlations between variables.

(Also not good!)

$$\Sigma = \sigma^2 I = \begin{pmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{pmatrix}$$

(Even harsher constraint!)

$$z \sim N(0, I) \quad z \in \mathbb{R}^d \quad (\text{choose } d < n)$$

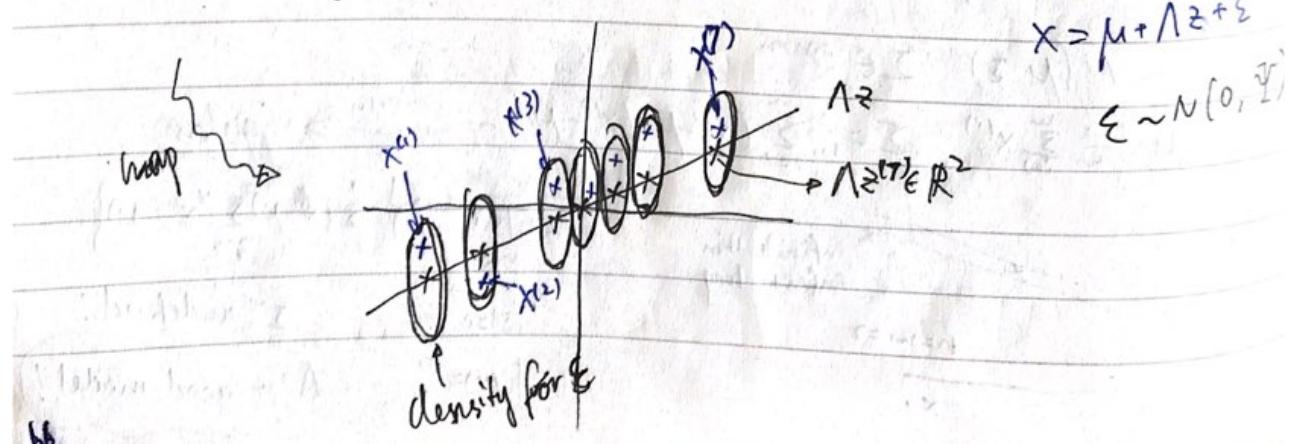
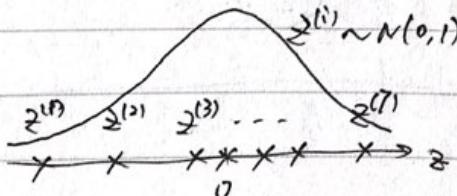
$$x|z \sim N(\mu + \Lambda z, \Psi)$$

Equivalently,

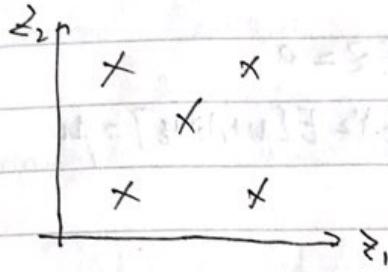
$$x = \mu + \Lambda z + \varepsilon, \quad \text{where } \varepsilon \sim N(0, \Psi)$$

Parameters: $\mu \in \mathbb{R}^n$, $\Lambda \in \mathbb{R}^{n \times d}$, $\Psi \in \mathbb{R}^{n \times n}$ (diagonal)

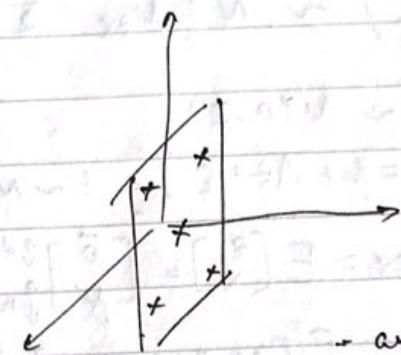
e.g.: $z \in \mathbb{R}^1$, $x \in \mathbb{R}^2$ $\Lambda = \begin{bmatrix} 2 & \\ & 1 \end{bmatrix}$ $\Psi = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$ $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$



e.g2: $z \in \mathbb{R}^2, x \in \mathbb{R}^3$



map
map

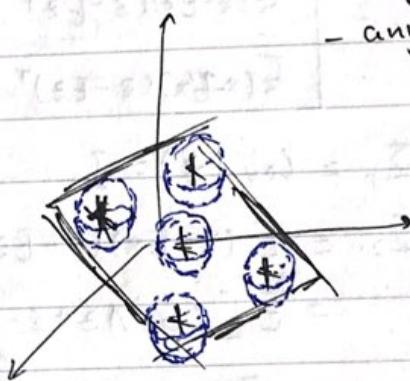


- any position
- any orientation

$$\mu + \Lambda z$$

$$x = \mu + \Lambda z + \varepsilon$$

\uparrow
3D Gaussian parallel to axis
(sphere)



If $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ $x_1 \in \mathbb{R}^r$ $x_2 \in \mathbb{R}^s$
 $x \in \mathbb{R}^{r+s}$

$$x \sim N(\mu, \Sigma) \text{ where } \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \begin{matrix} \overset{r}{\uparrow} & \overset{s}{\uparrow} \\ \downarrow & \downarrow \end{matrix} \quad \text{E.g.: } \Sigma_{12} \in \mathbb{R}^{rxs}$$

$$P(x)$$

$$P(x_1) = \int_{x_2} P(x_1, x_2) dx_2$$

$$x_1 \sim N(\mu_1, \Sigma_{11})$$

$$P(x_1 | x_2) = \frac{P(x_1, x_2)}{P(x_2)} \leftarrow N(\mu_1, \Sigma_{11})$$

$$x_1 | x_2 \sim N(\mu_{1|2}, \Sigma_{1|2})$$

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu_2)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

$$\begin{pmatrix} z \\ x \end{pmatrix} \sim N(\mu_{zx}, \Sigma)$$

$$z \sim N(0, I)$$

$$Ez = 0$$

$$x = \mu + \Lambda z + \varepsilon, \quad \varepsilon \sim N(0, \Sigma) \quad Ex = E[\mu + \Lambda z + \varepsilon] = \mu$$

$$\mu_{zx} = E\begin{pmatrix} z \\ x \end{pmatrix} = \begin{bmatrix} \vec{0} \\ \mu \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} E(z - Ez)(z - Ez)^T & E(z - Ez)(x - Ex)^T \\ E(x - Ex)(z - Ez)^T & E(x - Ex)(x - Ex)^T \end{bmatrix}$$

$$\Sigma_{11} = \text{Cov}(z) = I$$

$$\Sigma_{21} = E(x - Ex)(z - Ez)^T$$

$$= E[(\mu + \Lambda z + \varepsilon - \mu) \cdot z^T]$$

$$= E[\Lambda z z^T] - E[\varepsilon z^T]$$

$$= \Lambda E[z z^T] = \Lambda \quad \because \varepsilon, z \text{ independent.}$$

$$\Sigma_{22} = E[(\mu + \Lambda z + \varepsilon - \mu)(\mu + \Lambda z + \varepsilon - \mu)^T]$$

= ...

$$= \Lambda \Lambda^T + \Sigma$$

$$\text{So, } \begin{pmatrix} z \\ x \end{pmatrix} \sim N\left(\begin{bmatrix} \vec{0} \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Sigma \end{bmatrix}\right)$$

$$\{x^{(1)}, \dots, x^{(n)}\}$$

$$l(\theta) = \sum_i \log P(x^{(i)})$$

$$P(x^{(i)}) \quad x^{(i)} \sim N(\mu, \Lambda \Lambda^T + \Sigma)$$

E-step:

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

M-step:

$$\theta := \arg \max_{\theta} \sum_i \int_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} dz^{(i)}$$

Lecture 14

Factor Analysis

- EM steps

Principal Component Analysis (PCA)

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \quad X \sim N(\mu, \Sigma) \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

$$\text{Marginal: } P(X_1) = ? \quad X_1 \sim N(\mu_1, \Sigma_{11})$$

$$\text{Conditional: } P(X_1 | X_2) = ? \quad X_1 | X_2 \sim N(\mu_{12}, \Sigma_{12})$$

$$\mu_{12} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (X_2 - \mu_2)$$

$$\Sigma_{12} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

$$\{X^{(1)}, \dots, X^{(m)}\} \quad X^{(i)} \in \mathbb{R}^n \quad \text{Want } P(x)$$

$$z \sim N(0, I) \quad z \in \mathbb{R}^d$$

$$x = \mu + \Lambda z + \varepsilon \quad \varepsilon \sim N(0, \Sigma) \quad x | z \sim N(\mu + \Lambda z, \Sigma)$$

Parameters: $\mu \in \mathbb{R}^n$, $\Lambda \in \mathbb{R}^{n \times d}$, $\Sigma \in \mathbb{R}^{n \times n}$, diagonal

$$\left(\begin{array}{c} z \\ x \end{array} \right) \sim N(\mu_{zx}, \Sigma)$$

$$\left(\begin{array}{c} z \\ x \end{array} \right) \sim N\left(\begin{bmatrix} \bar{\mu} \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ \Lambda & \Lambda \Lambda^T + \Sigma \end{bmatrix} \right)$$

$$P(x) = ? \quad \checkmark \quad \checkmark$$

$$x \sim N(\mu, \Lambda \Lambda^T + \Sigma)$$

$$L(\theta) = \prod_{i=1}^m P(x^{(i)}; \mu, \Sigma, \Lambda) = \prod_{i=1}^m \frac{1}{(2\pi)^{\frac{n}{2}} |\Lambda^T + \Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu)^T (\Lambda^T + \Sigma)^{-1} (x^{(i)} - \mu) \right\}$$

If maximize log-likelihood (take derivative and set to 0), cannot get analytic solution. Instead, use EM algorithm.

E-step:

$$\text{Find } Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta)$$

M-step:

$$\theta := \arg \max_{\theta} \sum_i \int_{Z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} dz^{(i)}$$

① $z^{(i)}$ continuous variable.

$Q_i(z^{(i)})$ continuous PDF.

E-step (Computes $Q_i(z^{(i)})$):

$$z^{(i)}|x^{(i)}; \theta \sim N(\mu_{z^{(i)}|x^{(i)}}, \Sigma_{z^{(i)}|x^{(i)}}).$$

where

$$\mu_{z^{(i)}|x^{(i)}} = \bar{x} - \Lambda^T (\Lambda \Lambda^T + \Sigma)^{-1} (x^{(i)} - \mu)$$

$$\text{“} \mu_1 - \Sigma_{12} \Sigma_{22}^{-1} (x_2 - \mu) \text{”}$$

$$\Sigma_{z^{(i)}|x^{(i)}} = I - \Lambda^T (\Lambda \Lambda^T + \Sigma)^{-1} \Lambda$$

$$\text{“} \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \text{”}$$

M-step:

② $\int_{\mathbb{R}^n} Q_i(z^{(i)}) z^{(i)} dz^{(i)} \rightarrow \text{unnecessary complicated:}$

$$\int_{\mathbb{R}^n} \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_{z^{(i)}|x^{(i)}}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (z^{(i)} - \mu_{z^{(i)}|x^{(i)}})^T \Sigma_{z^{(i)}|x^{(i)}} (z^{(i)} - \mu_{z^{(i)}|x^{(i)}}) \right\} dz^{(i)}$$

simpler:

$$E_{z^{(i)} \sim Q_i} [z^{(i)}] = \mu_{z^{(i)}|x^{(i)}}$$

$$E_{z^{(i)} \sim Q_i} \left[\log \frac{P(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})} \right]$$

$$= E_{z^{(i)} \sim Q_i} \left[\log P(x^{(i)}|z^{(i)}; \theta) \right] \quad \text{focus on this (only one depends on parameters)}$$

$$+ E_{z^{(i)} \sim Q_i} \left[\log \frac{P(z^{(i)})}{Q_i(z^{(i)})} \right] \quad \text{many have parameter in other case}$$

will never have parameter
(it is fixed in the E-step)

e.g.: $\max_{\Lambda} \sum_{i=1}^m E_{z^{(i)}} \left[\log P(x^{(i)}|z^{(i)}; \Lambda, \Sigma, \mu) \right]$

$$x^{(i)}|z^{(i)} \sim N(\mu + \Lambda z^{(i)}, \Sigma) \quad \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x^{(i)} - \mu - \Lambda z^{(i)})^T \Sigma^{-1} (x^{(i)} - \mu - \Lambda z^{(i)}) \right\}$$

$$= \sum_{i=1}^m E \left[\text{constant} - \frac{1}{2} (x^{(i)} - \mu - \Lambda z^{(i)})^T \Sigma^{-1} (x^{(i)} - \mu - \Lambda z^{(i)}) \right]$$

$\partial_{\Lambda} (-\dots) \underset{\text{set}}{=} 0:$

$$\textcircled{3} \quad \Lambda = \left(\sum_{i=1}^m (x^{(i)} - \mu) \underbrace{E[z^{(i)\top}]}_{\text{“} \Sigma \text{”}} \left(\sum_{i=1}^m E[z^{(i)} z^{(i)\top}] \right)^{-1} \right)$$

$$\textcircled{3} \quad \Lambda = \left(\sum_{i=1}^m (x^{(i)} - \mu) E[z^{(i)\top}] \right) \left(\sum_{i=1}^m E[z^{(i)} z^{(i)\top}] \right)^{-1}$$

Expectations are w.r.t. $z^{(i)} \sim Q_i$

$$E[z^{(i)\top}] = \mu_{z^{(i)\top}|x^{(i)}}^T$$

$$\text{If } z^{(i)} \sim N(\mu, \Sigma) - \Sigma = Ezz^\top - EzEz^\top$$

$$\Rightarrow Ezz^\top = \Sigma + EzEz^\top$$

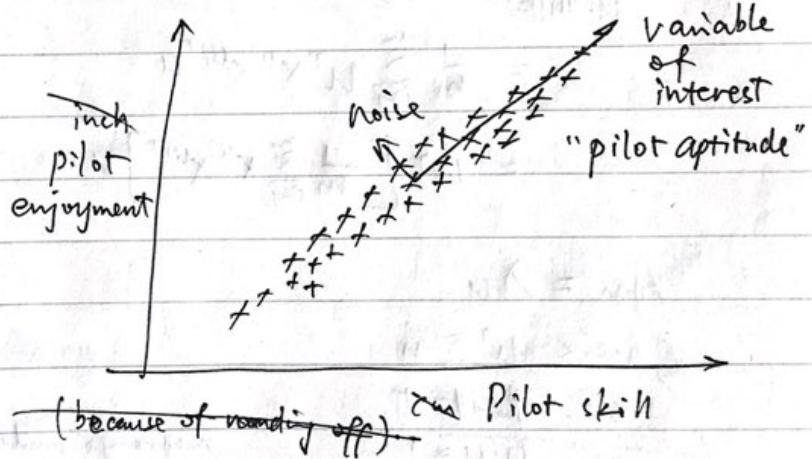
$$E[z^{(i)\top} z^{(i)\top}] = \Sigma_{z^{(i)\top}|x^{(i)}} + \mu_{z^{(i)\top}|x^{(i)}} \mu_{z^{(i)\top}|x^{(i)}}^T$$

Go to lecture notes
for skipped derivation

Principal Component Analysis (PCA)

Given $\{x^{(1)}, \dots, x^{(m)}\} \quad x^{(i)} \in \mathbb{R}^n$

Reduce it to k-dim data ($k < n$, often $k \ll n$)



Preprocessing

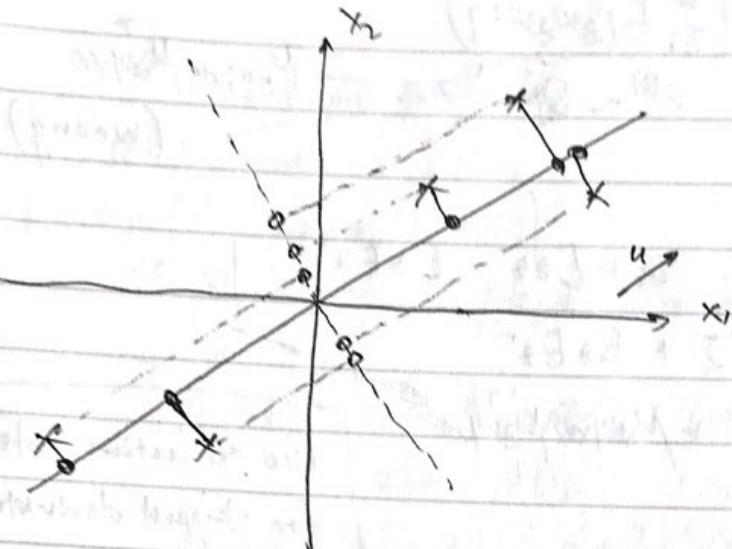
1. Set $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$ } zero out

2. Replace $x^{(i)}$ with $x^{(i)} - \mu$ } mean

3. Set $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)})^2$ } Normalize to unit variance

4. Replace $x_j^{(i)}$ with $\frac{x_j^{(i)}}{\sigma_j}$

Only do this step when
"different feature scales"



Objective:

Find a direction,
project data onto it,
get large variance.

If $\|u\|=1$, vector $x^{(i)}$ projected onto u has length $x^{(i)\top} u$

Choose u :

$$\max_{\|u\|=1} \frac{1}{m} \sum_{i=1}^m (x^{(i)\top} u)^2$$

$$= \frac{1}{m} \sum_{i=1}^m u^\top x^{(i)} x^{(i)\top} u$$

$$= u^\top \left[\frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)\top} \right] u \quad \Rightarrow \text{u is the principal eigenvector of } \Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)\top}$$

$$Au = \lambda u$$

$$\max_u u^\top \Sigma u$$

$$\text{st. } \|u\|=1$$

$$u^\top u = 1$$

Lagrange multiplier

$$\mathcal{L}(u, \lambda) = u^\top \Sigma u - \lambda(u^\top u - 1)$$

$$\nabla_u \mathcal{L} = \Sigma u - \lambda u \stackrel{!}{=} 0;$$

If want k -dim subspace,

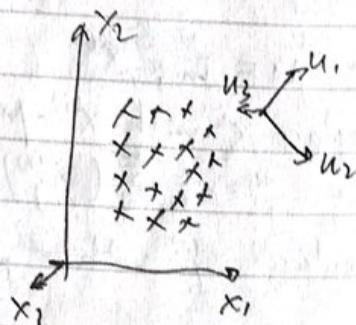
Choose u_1, \dots, u_k to be k

top eigenvectors of Σ .



Repeated eigenvalues:

eigenvectors can rotate freely
within their subspaces. Sometimes
it is dangerous to look one
basis at a time.



$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \alpha_0 \end{bmatrix}$$

