Q.

a) No. A discriminative classifier does not depend on any priors or probability distribution functions. It can only classify samples but not ~~generate~~ samples.
generate

b) Anomaly detection, because the number of failing engines are too small.

c) The prior can change the position of the decision boundary of LDA.

d)  1. ~~Normalize~~ Normalize the image data.
    2. Compute the covariance matrix of the data.
    3. Perform SVD on this matrix.
    4. Project the data onto top K eigen vectors to get the compression image

The K will control the degree of compression.

e)  $$\min \left\{ \frac{1}{2} \|w\|^2 + \lambda \sum_{i=1}^{n} \xi_i \right\}$$

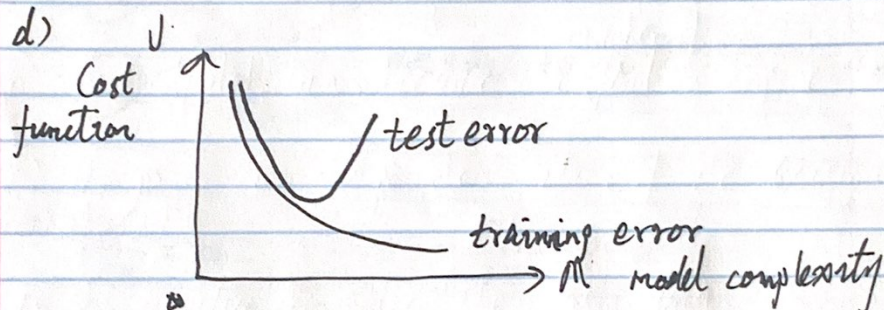    s.t. $(w^T x_i + b) y_i \geq 1 - \xi_i \ \forall i$
    $\xi_i \geq 0$

Q2

a) A: $\eta = 10^2$    B: $\eta = 10^0$    C: $\eta = 10^{-2}$

b) Model A has highest bias.

c) Model C has high variance.

d)



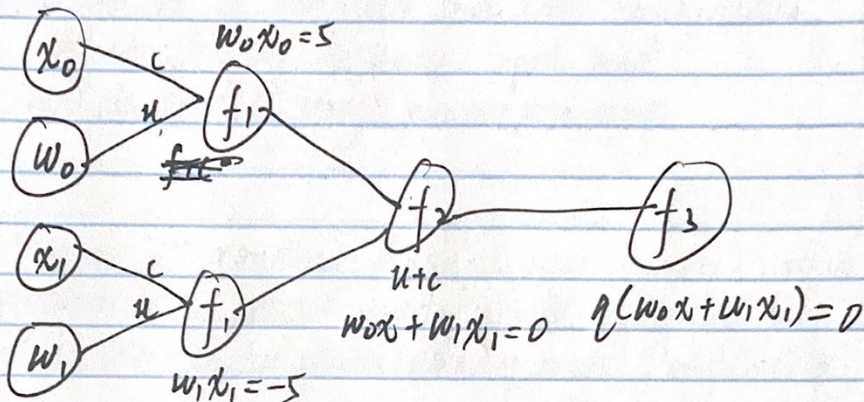Focus on the test error curve, when the curve goes up, overfitting is happening although trainning error is still going down.

Q3
a)



$W_0 x_0 = 5$

$u + c$
$W_0 x + W_1 x_1 = 0$   $q(W_0 x + W_1 x_1) = 0$

$W_1 x_1 = -5$

b)   $\frac{\partial f_1}{\partial u} = c$   $\frac{\partial f_2}{\partial u} = 1$   $\frac{\partial f_3}{\partial u} = 1 - u^2$

c)   $w x^T = |1 \times 5 - 5 \times 1| = 0$
$q(w x^T) = 0$   is the forward pass output.

d)   $h(x) = f(W_0 x + W_1 x_1)$
$h(x) = q(u)$   $u = W_0 x_0 + W_1 x_1$
$\frac{\partial h}{\partial W_0} = \frac{\partial h}{\partial u} \frac{\partial u}{\partial W_0} = [1 - (W_0 x_0 + W_1 x_1)^2] x_0$

$\frac{\partial h}{\partial x_0} = \frac{\partial h}{\partial u} \frac{\partial u}{\partial x_0} = [1 - (W_0 x_0 + W_1 x_1)^2] W_0$

$\frac{\partial h}{\partial W_1} = \frac{\partial h}{\partial u} \frac{\partial u}{\partial W_1} = [1 - (W_0 x_0 + W_1 x_1)^2] x_1$

$\frac{\partial h}{\partial x_1} = \frac{\partial h}{\partial u} \frac{\partial u}{\partial x_1} = [1 - (W_0 x_0 + W_1 x_1)^2] W_1$

Q4

a) It can be a ~~an~~ RNN and CNN combination.
RNN can handle sequence input and
CNN is good at image ~~pro do~~ processing.

b) Dropout. It randomly disables some neurons at a certain
probability to prevent overfitting. ~~When By~~ By doing that,
it is similar to ~~decre~~ reducing model complexity.

c) No.
Deeping learning itself takes care of feature engineering. ∧
Other machine learning algorithms require ~~not~~ manually.
                                                    feature engineering

d) No.
The function of activation is to ~~apply~~ make the model to
be able to learn non-linear features of the data. If you use
a linear activation function, it is hard for a model ~~the~~ to form
non-linear features by performing linear ∧ transformation.

e) Use sythetic images to augment the ~~image~~ data set,
such as cropping, coloring, rotation and ~~see~~ so on, which
can tell the model that these are also the features you should
learn.

f) CNN has covolutional layer and pooling layer.
CNN can better handle high-dimension-input like 2D image.

g)

Q5

a) $p(D|M) = \prod_{i=1}^{m} M^{x^{(i)}} (1-M)^{1-x^{(i)}}$

b) $p(M|D) = \dfrac{p(D|M)\, p(M)}{p(D)}$

$\underset{u}{\text{argmax}}\ p(M|D) = p(D|M)\, p(M)$

$=$

Q6

a) i) That means the prediction is correct.

ii) The prediction ~~is~~ lies on the decision boundary.

iii) The prediction is wrong.

b) ~~The loss function is not convex~~
   It is hard to find a minimum loss.

~~Cannot produce a reasonable solution.~~
~~It increases monotonically.~~
OR It can be a reasonable loss function,
   but not a good loss function,
   because we expect $z$ to be greater to classify ~~it~~ correctly.

c) i)

~~$J = --exp(-2z-1) + \frac{1}{2}\eta\|w\|^2$~~

$J = exp(-2z-1) + \frac{1}{2}\eta\|w\|^2$

ii) It should be minimized.

GD:
   Initialize $w$.
   Repeat:
$$\theta_{t+1} = \theta_t - \alpha\left[\frac{dJ}{dw} + \lambda\|w\|\right]$$

   Until it converges.