## Ex 2-1

a)

cost function:

$$\frac{1}{2}\sum_{i=1}^{N}(y_i - \hat{y}_i)^2$$

gradient descent - based method

**Ex 1-1:**

a) $\vec{1}^T \implies \vec{1}$

b) $2\vec{x}^T \to 2\vec{x}$

c) $2(\vec{x}-\vec{\mu})^T \to 2(\vec{x}-\vec{\mu})$

$$\vec{w} \leftarrow \vec{w} - \eta \cdot \frac{\partial\, cost}{\partial \vec{w}} = \vec{w} - \eta \cdot \frac{\partial}{\partial\vec{w}}\left[\sum_{i=1}^{N}\frac{1}{2}(y_i - \hat{y}_i)^2\right]$$

$$= \vec{w} - \eta \cdot \sum_{i=1}^{N}\frac{\partial}{\partial\vec{w}}\left[\frac{1}{2}(y_i - \hat{y}_i)^2\right]$$

$$= \vec{w} - \eta \cdot \sum_{i=1}^{N}(y_i - \hat{y}_i)\cdot\frac{\partial \hat{y}_i}{\partial\vec{w}} = \vec{w} - \eta \cdot \sum_{i=1}^{N}(y_i - \hat{y}_i)\cdot\vec{x}_i$$

Thus: $\quad \vec{w} \leftarrow \vec{w} - \eta \sum_{i=1}^{N}\left[(y_i - \hat{y}_i)\vec{x}_i\right] \quad \leftarrow$ a batch of sample.

Perceptron learning rule: $\quad \vec{w} = \vec{w} - \eta \cdot y_i \cdot \vec{x}_i \quad$ ( delta - rule )

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\uparrow \qquad\qquad$ ( sample - based )

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ one sample $\qquad\qquad\qquad$ rule

b) sample - based rule for gradient-based learning rule:

$$\vec{w} \leftarrow \vec{w} - \eta \cdot (y_i - \hat{y}_i)\vec{x}_i \quad \leftarrow \text{ stochastic gradient descent}$$

c) SGD can be learned on the fly. The model can be update sample by sample. It's unessary to recompute the whole model. Essentially better for large datased.
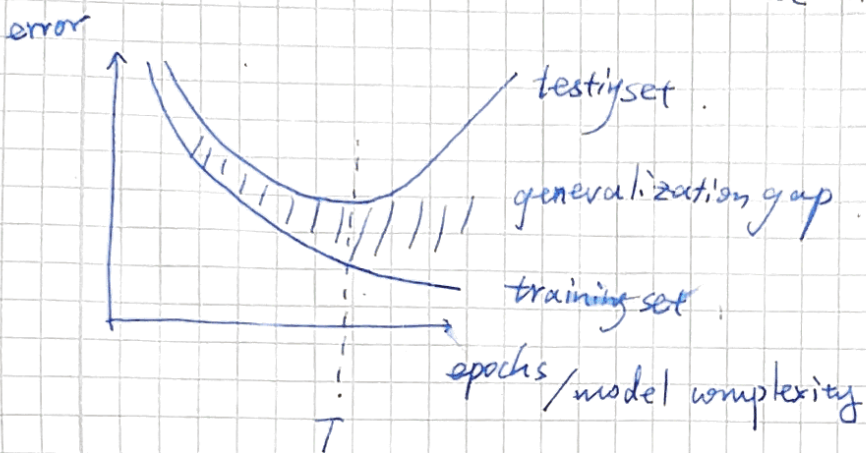
d) Striking difference : objective function .

- perceptron : binary outputs ⟶ classification
- ADALINE : real outputs ⟶ regression .
- If data is separable, perceptron converge faster than ADALINE .

⟶ error / loss / cost / objective function

Ex 2-2 .

a) overfitting : trained model over adapt to a given dataset .

Reason : ① generalization reason ( too many features )

② bias - variance decomposition



error

test set .
generalization gap
training set .
epochs / model complexity .

① when
# features >> # samples .
the model over adapts
to the sample while
training .

Bias - Variation decomposition:



error
variance
bias .

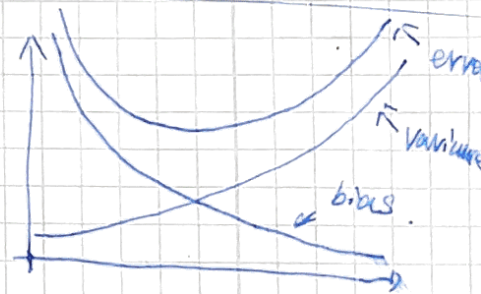$$error = \mathbb{E}[(y-\hat{y})^2] \quad (\text{mean square error})$$

$$= \mathbb{E}[(y - \mathbb{E}[\hat{y}] + \mathbb{E}[\hat{y}] - \hat{y})^2]$$

$$= \mathbb{E}[(y - \mathbb{E}[\hat{y}])^2] + \mathbb{E}[(\mathbb{E}[\hat{y}] - \hat{y})^2] + \mathbb{E}[2(y - \mathbb{E}[\hat{y}])(\mathbb{E}[\hat{y}] - \hat{y})]$$

ground truth   prediction   exp. of $\hat{y}$ prediction

$$= (y - \mathbb{E}[\hat{y}])^2 + \mathbb{E}[(\mathbb{E}[\hat{y}] - \hat{y})^2] + 2(y - \mathbb{E}(\hat{y}))(\mathbb{E}[\hat{y}] - \mathbb{E}[\hat{y}])$$

$(bias)^2 + $ variance

b) ① very good prediction, but very bad on test dataset.

② model complexity measure : (typically, in traditional machine learning is the # parameters .

A model has high complexity may cause a overfitting problem . The generalization gap :

$$\text{Test error} \leq \text{training error} + \overset{\text{model}}{\text{complexity penalty}}$$

$$= O(\sqrt{\frac{\#\text{parameter}}{\#\text{samples}}})$$

$\rightsquigarrow$ Testerror $-$ Training error $<<$ model complexity term

c) $-$ Regularization . (weight decay)

$-$ ~~#features~~ $<<$ # samples

Ex 2-3

a) continues / differentiable almost every where .

b) $f(\vec{x_i}) = \sum\limits_{h=0}^{M_p-1} W_h \, \phi_h (\vec{x_i}, \vec{V_h}) = \sum\limits_{h=0}^{M_p-1} \left[ \sum\limits_{j=0}^{M} W_{h,j} \, x_{i,j} \right] \cdot W_h$

$= \sum\limits_{j=0}^{M} \left[ \sum\limits_{h=0}^{M_p-1} W_{h,j} \cdot W_h \right] \cdot x_{i,j} = \sum\limits_{j=0}^{M} \alpha_j \cdot x_{i,j}$ .

which is equivalent to linear model .

c) No . Radial basis function (RBF) : $\phi(x_i, , \mu_j, , \sigma_j) = exp \{\frac{\|x_{ij} - \mu_{jl}}{2\sigma_j^2}$

which are typically used as a first layer in a 2-layer RBF network. The second layer is a linear combination.
It is similar to a gaussian kernel density estimator.