

CSE 142 Final Exam Information Sheet

Manhattan (L1) distance: $d(x, y) = \sum_{i=1}^d |x_i - y_i|$

Euclidian (L2) distance: $d(x, y) = \|x - y\| = \left(\sum_{i=1}^d (x_i - y_i)^2 \right)^{1/2}$

Minkowski (Lp) distance: $d(x, y) = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{1/p}$

Chebyshev distance:

$$L_{\infty}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_{\infty} = \max_i |x_i - y_i|$$

Hamming distance:

$$L_0(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_0 = \text{count}(|x_i - y_i| > 0)$$

Mahalanobis distance:

$$D_M(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \mathbf{y})}$$

Sample mean: $\hat{\mu}_x = \frac{1}{n} \sum_i x_i$

Sample variance: $\hat{\sigma}_x^2 = \frac{1}{n} \sum_i (x_i - \hat{\mu}_x)^2$

Sample covariance: $\hat{\sigma}_{xy} = \frac{1}{n} \sum_i (x_i - \hat{\mu}_x)(y_i - \hat{\mu}_y)$

Sample covariance matrix: $\hat{\Sigma} = \frac{1}{k} X_z X_z^T = \frac{1}{k} S$ where S is the scatter matrix

Gram matrix
(X is not zero-centered) $G = \mathbf{X}^T \mathbf{X}$

Bayes Rule:

$$P(H_i | D) = \frac{P(D | H_i) P(H_i)}{P(D)}$$

$$\text{False positive rate (FPR)} = \frac{FP}{N} = \alpha$$

$$\text{Accuracy} = \frac{TP + TN}{P + N} = \left(\frac{P}{P + N}\right) TPR + \left(\frac{N}{P + N}\right) TNR$$

$$\text{False negative rate (FNR)} = \frac{FN}{P} = \beta$$

$$\text{Error rate} = \frac{FP + FN}{P + N}$$

$$\text{True positive rate (TPR)} = \frac{TP}{P} = \text{Recall}$$

$$\text{Precision} = \frac{TP}{\hat{P}}$$

$$\text{True negative rate (TNR)} = \frac{TN}{N}$$

$$\text{Accuracy} + \text{error rate} = 1$$

$$F1 \text{ score} = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2 \cdot TP}{P + \hat{P}}$$

Scoring classifier margin: $z(x) = c(x) \hat{s}(x)$ (true class function * score)

Margin loss function: $L(z(x)) \rightarrow [0, \infty)$

Ranking classifier error rate: $\text{rank-err} = \text{err} / PN$

Ranking classifier accuracy: $\text{rank-acc} = 1 - \text{rank-err}$

$$\text{Laplace correction} = \frac{N_i + 1}{|S| + k}$$

$$\text{m-estimate} = \frac{N_i + m\pi_i}{|S| + m}$$

Size of hypothesis space: $|H| = 2^{(\# \text{ instances})}$

Ranking classifier error rate: $\text{rank-err} = \text{err} / pN$

Ranking classifier accuracy: $\text{rank-acc} = 1 - \text{rank-err}$

Min. training set size for PAC learning: $m \geq \frac{1}{\varepsilon} \left(\ln|H| + \ln \frac{1}{\delta} \right)$

PAC learning outputs, with probability at least $1 - \delta$, a hypothesis h such that $\text{err}_D < \varepsilon$

Multivariate least-squares regression
(homogeneous representation)

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \mathbf{S}^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Least-squares minimization with regularization:

$$\mathbf{w}^* = \text{argmin} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda r(\mathbf{w})$$

Impurity measures:

Minority class

$$\text{Imp}(\dot{p}) = \min(\dot{p}, 1 - \dot{p})$$

Gini index

$$\text{Imp}(\dot{p}) = 2\dot{p}(1 - \dot{p})$$

Entropy

$$\text{Imp}(\dot{p}) = -\dot{p} \log_2(\dot{p}) - (1 - \dot{p}) \log_2(1 - \dot{p})$$

$\sqrt{\text{Gini index}}$

$$\text{Imp}(\dot{p}) = \sqrt{2\dot{p}(1 - \dot{p})}$$

Total impurity:

$$\text{Imp}(\{D_1, \dots, D_l\}) = \sum_{i=1}^l \frac{|D_i|}{|D|} \text{Imp}(D_i)$$

Proportion of positive instances
in a (binary) data partition:

$$\dot{p} = \frac{P}{P + N}$$

In a k-class data partition:

$$\dot{p}_i = \frac{C_i}{\sum_{i=1}^k C_i}$$

Classifier margin for point \mathbf{x}

$$\mathbf{z}(\mathbf{x}) = \frac{y(\mathbf{w}^T \mathbf{x} - t)}{\|\mathbf{w}\|} = \frac{m}{\|\mathbf{w}\|} \quad \text{Non-homogeneous representation}$$

Algorithm $\text{GrowTree}(D, F)$ – grow a feature tree from training data.

Input : data D ; set of features F .

Output : feature tree T with labelled leaves.

if $\text{Homogeneous}(D)$ **then return** $\text{Label}(D)$;

$S \leftarrow \text{BestSplit}(D, F)$; // e.g., BestSplit-Class (Algorithm 5.2)

split D into subsets D_i according to the literals in S ;

for each i **do**

if $D_i \neq \emptyset$ **then** $T_i \leftarrow \text{GrowTree}(D_i, F)$;

else T_i is a leaf labelled with $\text{Label}(D)$;

end

return a tree whose root is labelled with S and whose children are T_i

Algorithm $\text{Perceptron}(D, \eta)$ – train a perceptron for linear classification.

Input : labelled training data D in homogeneous coordinates; learning rate η .

Output : weight vector \mathbf{w} defining classifier $\hat{y} = \text{sign}(\mathbf{w} \cdot \mathbf{x})$.

$\mathbf{w} \leftarrow \mathbf{0}$; // Other initialisations of the weight vector are possible

$\text{converged} \leftarrow \text{false}$;

while $\text{converged} = \text{false}$ **do**

$\text{converged} \leftarrow \text{true}$;

for $i = 1$ to $|D|$ **do**

if $y_i \mathbf{w} \cdot \mathbf{x}_i \leq 0$ // i.e., $\hat{y}_i \neq y_i$

then

$\mathbf{w} \leftarrow \mathbf{w} + \eta y_i \mathbf{x}_i$;

$\text{converged} \leftarrow \text{false}$; // We changed \mathbf{w} so haven't converged yet

end

end

end

Soft margin optimization problem:

$$\mathbf{w}^*, t^*, \xi_i^* = \underset{\mathbf{w}, t, \xi_i}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i - t) \geq 1 - \xi_i$ and $\xi_i \geq 0, 1 \leq i \leq n$

Boosting

Confidence factor

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}$$

Misclassified
point

$$w' = \frac{w}{2\epsilon_t}$$

Correctly classified
point

$$w' = \frac{w}{2(1 - \epsilon_t)}$$

Sigmoid function:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Backpropagation error for output units:

$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$

Backpropagation error for hidden units:

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{kh} \delta_k$$

Algorithm $K\text{Means}(D, K)$ – K -means clustering using Euclidean distance Dis_2 .

Input : data $D \subseteq \mathbb{R}^d$; number of clusters $K \in \mathbb{N}$.

Output : K cluster means $\mu_1, \dots, \mu_K \in \mathbb{R}^d$.

randomly initialise K vectors $\mu_1, \dots, \mu_K \in \mathbb{R}^d$;

repeat

assign each $\mathbf{x} \in D$ to $\operatorname{argmin}_j \text{Dis}_2(\mathbf{x}, \mu_j)$;

for $j = 1$ to K **do**

$D_j \leftarrow \{\mathbf{x} \in D \mid \mathbf{x} \text{ assigned to cluster } j\}$;

$\mu_j = \frac{1}{|D_j|} \sum_{\mathbf{x} \in D_j} \mathbf{x}$;

end

until no change in μ_1, \dots, μ_K ;

return μ_1, \dots, μ_K ;

Algorithm Bagging(D, T, \mathcal{A}) – train an ensemble of models from bootstrap samples.

Input : data set D ; ensemble size T ; learning algorithm \mathcal{A} .

Output : ensemble of models whose predictions are to be combined by voting or averaging.

```
for  $t = 1$  to  $T$  do
    build a bootstrap sample  $D_t$  from  $D$  by sampling  $|D|$  data points with
    replacement;
    run  $\mathcal{A}$  on  $D_t$  to produce a model  $M_t$ ;
end
return  $\{M_t | 1 \leq t \leq T\}$ 
```

Algorithm Boosting(D, T, \mathcal{A}) – train an ensemble of binary classifiers from reweighted training sets.

Input : data set D ; ensemble size T ; learning algorithm \mathcal{A} .

Output : weighted ensemble of models.

$w_{1i} \leftarrow 1/|D|$ for all $x_i \in D$; // start with uniform weights

```
for  $t = 1$  to  $T$  do
    run  $\mathcal{A}$  on  $D$  with weights  $w_{ti}$  to produce a model  $M_t$ ;
    calculate weighted error  $\epsilon_t$ ;
    if  $\epsilon_t \geq 1/2$  then
        set  $T \leftarrow t - 1$  and break
    end
     $\alpha_t \leftarrow \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$ ; // confidence for this model
     $w_{(t+1)i} \leftarrow \frac{w_{ti}}{2\epsilon_t}$  for misclassified instances  $x_i \in D$ ; // increase weight
     $w_{(t+1)j} \leftarrow \frac{w_{tj}}{2(1-\epsilon_t)}$  for correctly classified instances  $x_j \in D$ ; // decrease
end
return  $M(x) = \sum_{t=1}^T \alpha_t M_t(x)$ 
```

The last line should be:

return $M(x) = \text{sign}(\sum_{t=1}^T \alpha_t M_t(x))$