

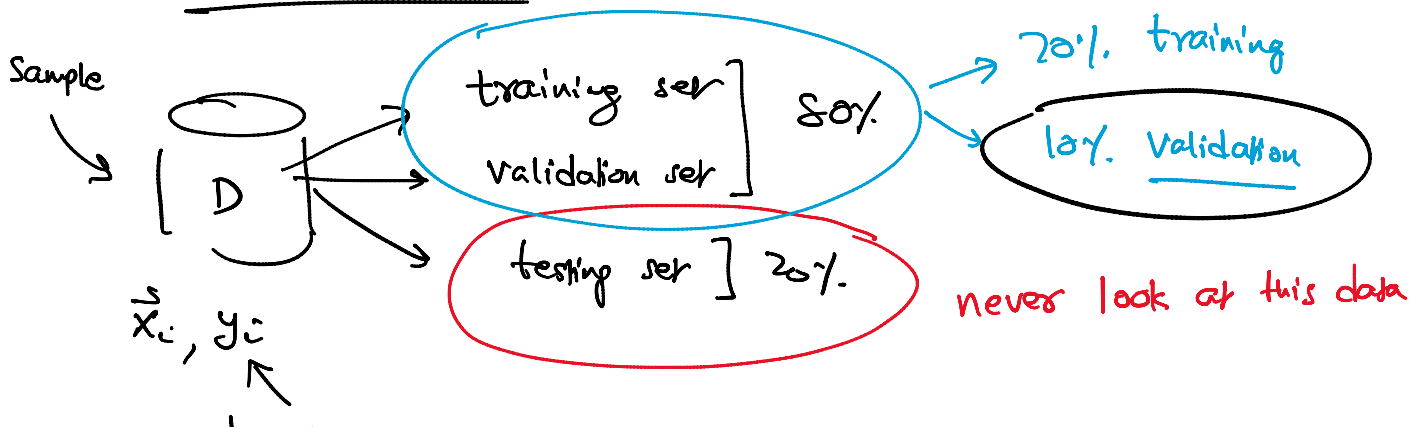
Lecture 17

Sunday, October 29, 2023 9:38 PM

Exam I (Thursday, Nov 2)

- 1) Logistic regression (Chap 24) — log odd ratio
 - 2) Neural Networks (Chap 25) — net gradients, ∇_w, ∇_b , output activations
 - 3) Bayes classifier (Chap 18) — naive, full (no categorical data)
 - KNN classifier
 - 4) Decision Trees (Chap 19) — split eval, metrics
 - 5) Support Vector Machines (Chap 21) —
 - a) margin vs. slack
 - b) given $\vec{\alpha}$, compute \vec{w} (linear)
 - polynomial (quadratic)
- $\vec{x} \rightarrow \underbrace{\phi(\vec{x})}_{\text{feature space}}$
- $\vec{w} = \sum_{i=1}^n \alpha_i y_i \phi(x_i)$
- c) some iterations of SGA (find $\vec{\alpha}$)
- 6) Classification assessment (see 22.1 only)
 - prec, recall, f1, Roc (TPR vs FPR)

Classification Assessment



..., y_c
true response

User-specified constants

↳ hyperparameters

↓

← C value in SVMs

learning rate

try diff C

values

(grid search)

$\{ 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4, \dots \}$

build model on training

eval on validation

Once we have a trained & validated model

↳ then apply on test data

↳ compute F1-score

k-fold cross validation (cv)

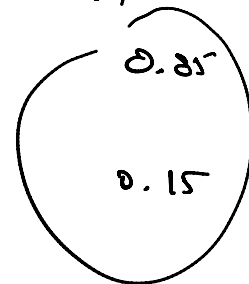
↳ repeat k times: build k different models

↳ compute Θ_i $i = 1 \dots k$ (on the test data)

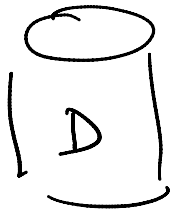
Some measure, e.g. accuracy, F1

↳ $\bar{\mu} = \text{mean}\{\Theta_1, \Theta_2, \dots, \Theta_k\}$

↳ $\frac{1}{k} \sum \Theta_i^2 = \text{Var}\{\Theta_1, \Theta_2, \dots, \Theta_k\}$



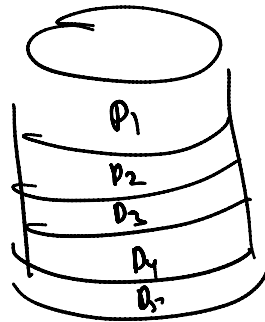
$$0.85 \pm \sqrt{0.15}$$



Sample

$k=5$
5-fold CV

- 1) shuffle D
- 2) create 5 equal partitions



5 - folds

- 3) keep aside D_i for testing
Use remaining folds for training & validation

$$D \setminus D_i = D - D_i$$

$$T_i = D \setminus D_i \leftarrow \text{training}$$

$$D_i \leftarrow \text{testing}$$

$$\theta_i \leftarrow \text{Model}(D_i)$$

test	train
D_1	$D_2 D_3 D_4 D_5$
D_2	$D_1 D_3 D_4 D_5$
\vdots	
D_5	$D_1 D_2 D_3 D_4$

LOOCV : leave one out CV

↳ usually used for data which is expensive to collect



n-fold CV

1 point for testing

n-1 points for training

1 point for testing \rightarrow $n-1$ points for training

→ $n-1$ points for training

Bootstrap sampling



Sample

Random dataset of size n

D:
$$\forall i=1 \dots k$$

→ Sample with replacement

each point can be selected multiple times

↳ training set: D_i (further split into validation)

↳ test on D

$$D_i \subseteq D$$

$\theta_i \leftarrow$ over-estimate (optimistic estimate)

Since we have seen many of the top points in D_i

Q: Sample with replacement, n trials

What is the probability that a point will not be sampled?

c) what is the prob that a point is sampled

$$\frac{1}{n}$$

b) prob that 'it' is not sampled

$$\left(1 - \frac{1}{n}\right) \text{ for one trial}$$

c) not sampled even after n trials

[illegible]

$$P(x_j \notin D_i) = \left(1 - \frac{1}{n}\right)^n \approx e^{-1} = 0.368$$

a point is
not sampled
for D_i

D_i contains only about 63.2% of the points!

From $\hat{\mu}$ & $\hat{\sigma}^2$ (mean & variance) for a classifier
estimates

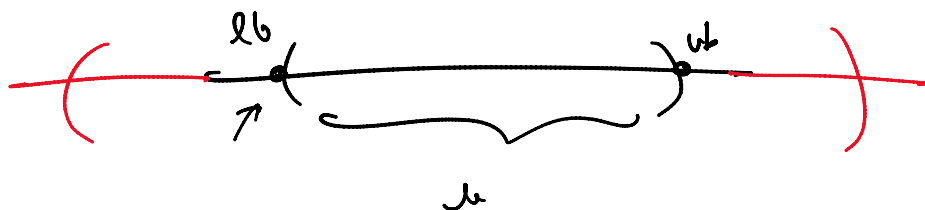
we get a confidence bound on the
true mean

expected performance (accuracy)
f1

μ
true mean

Compute the confidence interval for the true mean

$$\left(\underbrace{\hat{\mu} - t_k^\alpha(\hat{\sigma})}_{lb} \leq \mu \leq \underbrace{\hat{\mu} + t_k^\alpha(\hat{\sigma})}_{ub} \right)$$



α : confidence level
95%

t_k : student's t-distribution

α : confidence level
90%.

$\alpha = 99\%$

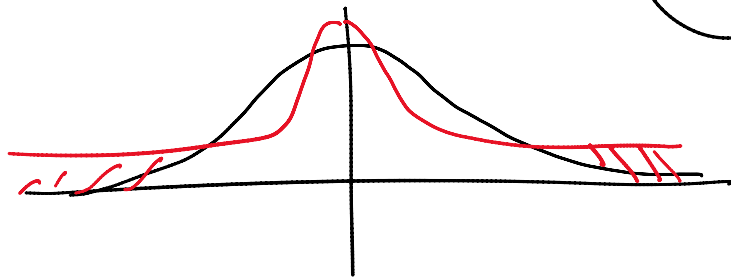
t_k : student's t-distribution
with $k-1$ degrees of freedom

t_k is like the normal
as $k \rightarrow \infty$

t_k is small sample version of normal

k is small e.g.

$k=5$
 $k=10$



Q How to compare two models?

M_A & M_B
(svm) (logistic regression)

$\hat{\mu}_A \pm \hat{\sigma}_A$ vs $\hat{\mu}_B \pm \hat{\sigma}_B$



Paired t-test

CV k-folds

for $i = 1 \dots k$

$T_i = D \setminus D_i \leftarrow$ training
... ..

$T_i = D \setminus D_i \leftarrow \text{training}$
 $D_i = i\text{-th fold} \leftarrow \text{testing}$
 train M_A & M_B on T_i
 Accuracy A, B $\left\{ \begin{array}{l} \theta_i^A = M_A(D_i) \\ \theta_i^B = M_B(D_i) \end{array} \right.$
 $\delta_i = \theta_i^A - \theta_i^B$: difference in performance

$\hat{\mu}_\delta \equiv \text{mean difference} = \text{mean}\{\delta_1, \delta_2, \dots, \delta_k\}$

$\hat{\sigma}_\delta^2 \equiv \text{variance of difference}$

Hypothesis testing

Q. Are the two models different in terms of θ (accuracy)

H_0 : null hypothesis

$\left[\begin{array}{l} \text{a) there is no difference between the two models} \\ \text{b) look for evidence to support or reject } H_0 \end{array} \right]$

H_a : alternative hypothesis

\rightarrow there is a difference

H_0 : there is no difference

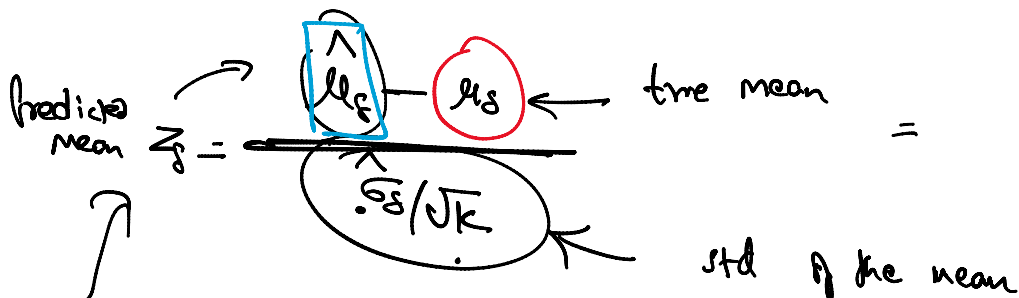
$\mu_\delta = 0$

true expected difference

$$\mu_S = 0$$

true expected difference

t-test:



Standardized difference (how many deviations away from true mean)

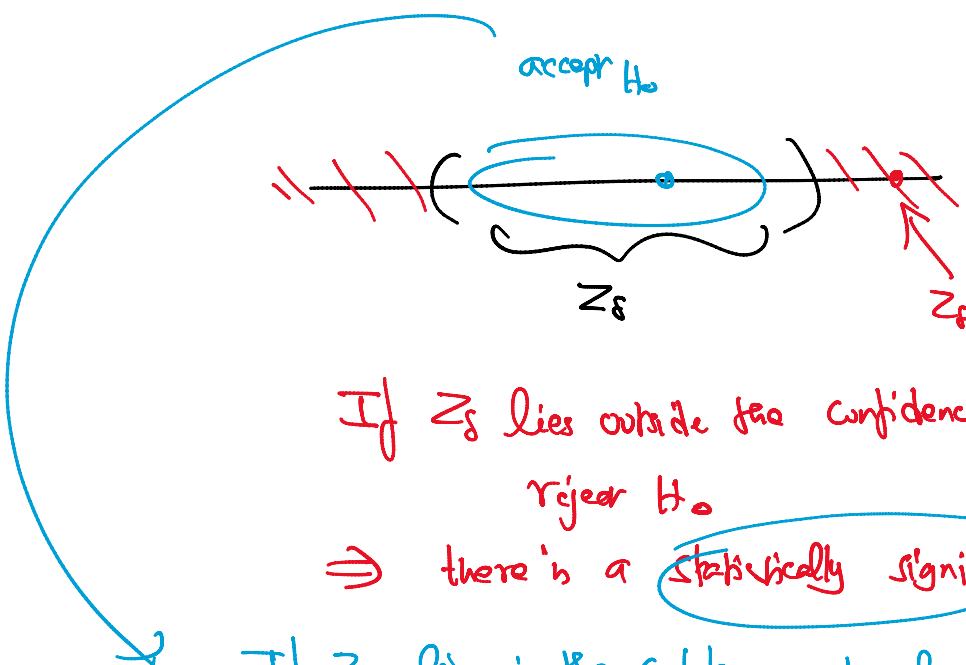
this has a t-distribution with $k-1$ degrees of freedom

$$\left(-t_{\alpha/2}^k \leq Z_S \leq +t_{\alpha/2}^k \right)$$

$\alpha: 90\%$

$\alpha: 99\%$

what is expected under the null hypothesis



If Z_S lies outside the confidence interval then reject H_0

\Rightarrow there is a statistically significant difference

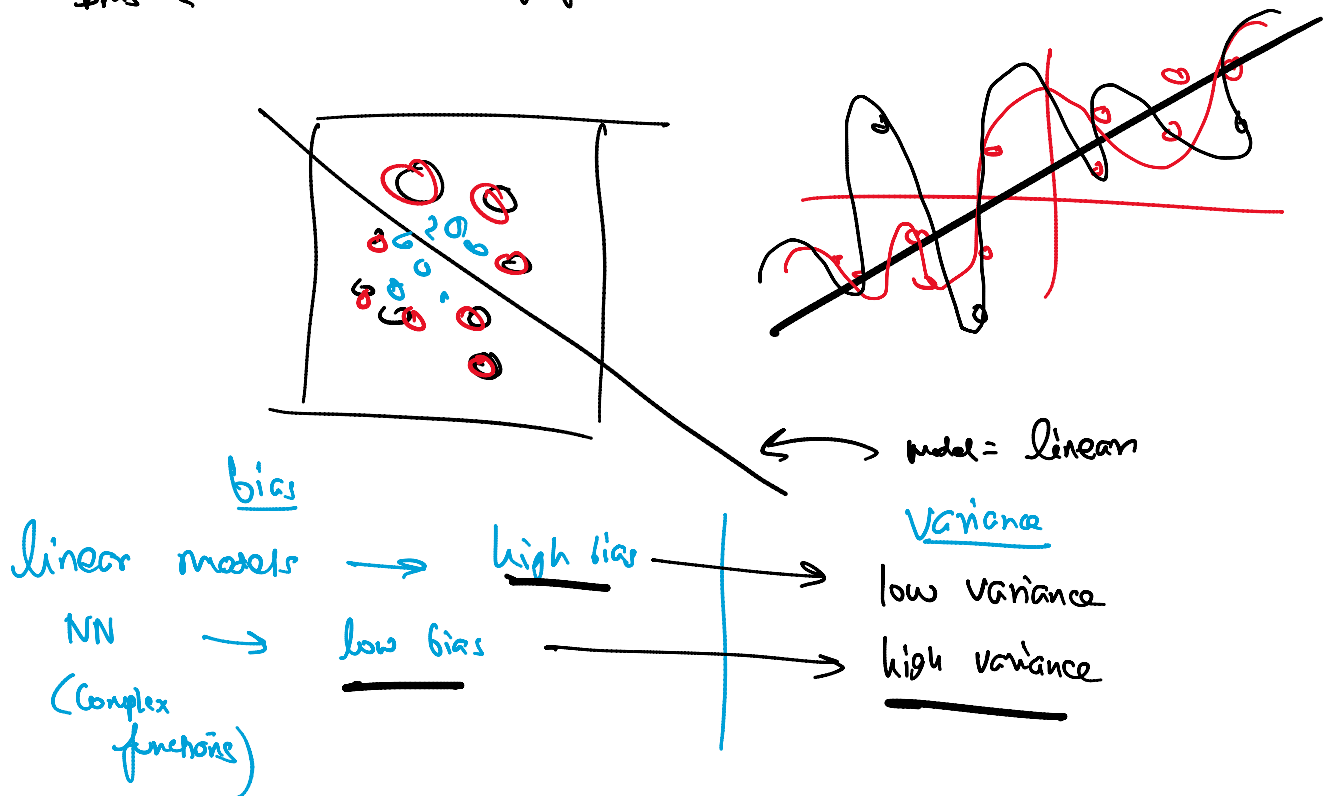
If Z_S lies in the confidence interval

$\Rightarrow H_0$ is probably the right conclusion

\Rightarrow It is probably the right conclusion
 \Rightarrow No difference!

Bias-Variance of the classifier

Bias \leftarrow Inherent class of functions a model can approximate.



$$\underline{\text{error}} = \text{bias} + \text{variance}$$

bias-variance tradeoff